# Foundations of Safe Online Reinforcement Learning in the Linear Quadratic Regulator: Generalized Baselines

Benjamin Schiffer and Lucas Janson

Department of Statistics, Harvard University

## Abstract

Many practical applications of online reinforcement learning require the satisfaction of safety constraints while learning about the unknown environment. In this work, we establish theoretical foundations for reinforcement learning with safety constraints by studying the canonical problem of Linear Quadratic Regulator learning with unknown dynamics, but with the additional constraint that the position must stay within a safe region for the entire trajectory with high probability. Our primary contribution is a general framework for studying stronger baselines of nonlinear controllers that are better suited for constrained problems than linear controllers. Due to the difficulty of analyzing non-linear controllers in a constrained problem, we focus on 1-dimensional state- and action- spaces, however we also discuss how we expect the high-level takeaways can generalize to higher dimensions. Using our framework, we show that for *any* non-linear baseline satisfying natural assumptions, $\tilde{O}_T(\sqrt{T})$-regret is possible when the noise distribution has sufficiently large support, and $\tilde{O}_T(T^{2/3})$-regret is possible for *any* subgaussian noise distribution. In proving these results, we introduce a new uncertainty estimation bound for nonlinear controls which shows that enforcing safety in the presence of sufficient noise can provide "free exploration" that compensates for the added cost of uncertainty in safety-constrained control.

1

# 1   Introduction

## 1.1   Background and Motivation

Recent advances in reinforcement learning (RL) have led to many successes in applying RL algorithms to a variety of practical online applications, from robotics to personalized health [Levine et al., 2016, Lillicrap et al., 2015, Tewari and Murphy, 2017]. A core concept behind online RL algorithms is the careful balance between exploration (proactively learning about the unknown environment) and exploitation (using what is already known to maximize reward). In practice, however, RL algorithms are restricted in the possible actions and states by safety constraints. For example, a drone using an RL algorithm must have safety constraints restricting possible states that would result in the drone crashing into a building or injuring a bystander. Therefore, the drone cannot explore the environment by accelerating directly into a building, and instead must explore in a safe manner. To deploy more RL algorithms to practical applications, instead of just balancing exploration versus exploitation, the optimal algorithm must now balance exploration versus exploitation versus safety. In many applications, the safety constraints must be obeyed at all time steps (even at the beginning), which does not allow for any violation of safety even during the initial learning period. Therefore, this component of "safety" involves both learning safely as well as learning how to be safe in the future. Studying simple canonical problems in RL can give insights into how to develop safe RL algorithms in more complex practical settings. In this paper, we address safety in the context of online LQR with unknown dynamics. Online LQR with unknown dynamics can be viewed as one of the simplest RL problems with a continuous decision space, and this problem has recently gained significant attention within the RL community both with and without safety constraints (see e.g. Abbasi-Yadkori and Szepesvári [2011], Dean et al. [2018, 2019]).

## 1.2   Setting and Motivation

In order to better understand the interaction between safety and the balance of exploration/exploitation, we study the classic problem of controlling a discrete-time linear dynamical system with unknown dynamics while minimizing a quadratic cost. In our problem setting, the position at the next time step depends on the current position, the current control input, and a random noise. The goal is to choose controls (actions) that keep the position as close to the origin as possible while using as little control as possible. An example application of this problem is controlling a drone around a target [Rubio et al., 2016]. In this scenario, the goal is to maintain a safe distance from the target while preserving fuel despite random disturbances from air currents. In this paper, we are interested in the setting where the dynamics are unknown. When the dynamics are unknown, LQR becomes an online RL problem of balancing exploration (controls that learn about the dynamics) and exploitation (controls that minimize the cost). Extending the previous example of controlling a drone around a target, the dynamics could for example be determined by the weather pattern that is unknown in advance. The goal in this paper is to design an algorithm that can learn the dynamics safely while not incurring significantly more cost than the best safe algorithm when the dynamics are known.

To quantify safety in this setting, we will consider constraints on the position of the controller which restrict the position to stay within a safe region. Continuing the previous example, a drone control must be safe in that it must avoid positions that are currently occupied by walls or other objects. We focus on position constraints rather than control constraints because position constraints have the added difficulty that, at the time of choosing the control, the next position for any given control is unknown due to the noise and uncertainty about the dynamics. In contrast, the algorithm has perfect information about (and control over) the choice of control. We therefore consider the LQR setting with only position constraints. See Section 5 for discussion on how our results extend to the setting with control constraints. While the optimal policy with known dynamics and without position constraints is the well-understood Linear Quadratic Regulator, with constraints the optimal policy even for known dynamics no longer has a closed-form [Rawlings and Mayne, 2009]. Due to the substantially increased complexity of the constrained LQR problem with both known and unknown dynamics, we will focus on the setting when both the positions and controls are one-dimensional. We focus on the one-dimensional setting to highlight the main ways in which learning unknown dynamics changes in the presence of constraints, without the additional technical overhead that comes with proving results for higher dimensions. However, we do predict that many of the results in this paper can be generalized to higher dimensions, and we discuss this further in Section 5. Other works have also taken the same approach of first studying only the one-dimensional case of LQR, see e.g. Fefferman et al. [2021], Abeille and Lazaric [2017]. The one-dimensional setting of safe LQR does have its own applications, for example maintaining a fixed temperature of a room [Oldewurtel et al., 2008]. In this application, the goal is to maintain a certain range of safe temperatures with high probability while using as little energy as possible. Taking temperature as the position, this problem can be formulated as a one-dimensional LQR problem with safety "position" constraints on the temperature.

## 1.3 Our Contribution

The main theorems of this paper each establish new regret results for safety-constrained LQR learning. We improve prior works' regret bounds for this setting along three dimensions, the regret rate, the regret baseline, and the types of noise distributions. In contrast to prior works, we focus on one-dimensional LQR with only positional constraints, but in this setting we study more general non-linear baselines. The following table summarizes our different results relative to prior works across these three dimensions:

|  | Regret Rate | Regret Baseline | Noise Distributions |
| --- | --- | --- | --- |
| Previous works | $\tilde{O}_T(T^{2/3})$ | Best Safe Linear Controller | Bounded |
| Theorem 1 | $\tilde{O}_T(\sqrt{T})$ | Best General Baseline Controller | subgaussian+Large Support |
| Theorem 2 | $\tilde{O}_T(T^{2/3})$ | Best General Baseline Controller | subgaussian |

The main contribution of this paper is a general framework for analyzing different baselines for the safety-constrained LQR learning problem. Using this framework, we show that a $\tilde{O}_T(\sqrt{T})$ rate of regret is possible for safety-constrained LQR learning in one-dimension

3

for noise distributions with large support, improving on $\tilde{O}_T(T^{2/3})$ regret of previous works [Li et al., 2021, Dean et al., 2019]. This rate of regret for constrained LQR learning matches the optimal regret rate for *unconstrained* LQR learning [Ziemann and Sandberg, 2024]. In addition to improving the rate of regret, this result is also with respect to a stronger and more general baseline than studied in previous works. The regret for this result is defined with respect to the best controller from general classes of baseline controllers satisfying only minimal regularity condition. To the best of our knowledge, this is the first work on constrained LQR learning with respect to any baseline other than the best safe linear controller. Our result also holds for any subgaussian noise distribution, which is the (to the best of our knowledge) first safety-constrained LQR learning result for unbounded distributions. A key technical tool used to prove $\tilde{O}_T(\sqrt{T})$ regret is a new bound for estimating unknown dynamics with non-linear controllers, which may be of independent interest.

In addition to showing that $\tilde{O}_T(\sqrt{T})$ regret is possible when the noise distribution has sufficiently large support, we also show that a certainty equivalence algorithm can achieve a regret rate of $\tilde{O}_T(T^{2/3})$ relative to the best controller from these general classes of baseline controllers for any subgaussian noise distribution. All of the proofs of the regret results in this paper are constructive and provide certainty equivalence algorithms for achieving the guaranteed rates of regret.

## 1.4 Related Work

RL has been recognized as being a powerful tool in a broad array of applications [Silver et al., 2016, Kiran et al., 2021, Levine et al., 2016], but there is still a need to better understand RL in the presence of safety constraints. There exists a wide array of definitions of safety in RL, many of which focus on some notion of reachability or stability [Ganai et al., 2024, Garg et al., 2024, Gu et al., 2022, Moldovan and Abbeel, 2012, Wachi et al., 2018, 2024, Yao et al., 2024]. However, these notions of safety are less directly related to our problem setting. More related to our problem, there is also a body of literature on algorithms for RL for control with constraints that maintain safety for the entire trajectory [Fulton and Platzer, 2018, Cheng et al., 2019, Marvi and Kiumarsi, 2021, Fisac et al., 2018]. These works study different broad definitions of safety in control, which can apply to a wider variety of models and settings than our results. However, the technical contribution of these works focuses specifically on developing safe algorithms, without proving theoretical results about the rates of regret or the optimality of the proposed safe algorithms.

The LQR problem has many applications despite the simplicity of the problem statement [Priess et al., 2014, Choi and Seo, 1999, Shabaani and Jalili-Kharaajoo, 2003]. There has recently been a large body of work focusing on minimizing regret in the unconstrained LQR setting with unknown dynamics, beginning with Abbasi-Yadkori and Szepesvári [2011] which gave the first algorithm for $\tilde{O}_T(\sqrt{T})$ regret for unconstrained LQR learning. This was followed by many works that study variations of both the infinite and finite time problem including (but not limited to) Dean et al. [2018], Mania et al. [2019, 2020], Simchowitz et al. [2018], Cohen et al. [2019], Wang and Janson [2021, 2022], Mania et al. [2019], Abeille and Lazaric [2017], Zheng and Li [2020], Sun et al. [2020], Khosravi and Smith [2020], Sattar and Oymak [2022], Faradonbeh et al. [2018a, 2017], Oymak and Ozay [2019], Ye et al. [2024], Athrey et al. [2024], Ziemann and Sandberg [2024], Lee et al. [2024]. Certainty

Equivalence (CE) algorithms estimate the unknown dynamics and find an optimal policy under the estimated dynamics. Later works on LQR learning showed that CE algorithms are in fact (rate) optimal for the unconstrained learning problem [Simchowitz and Foster, 2020, Faradonbeh et al., 2018b, Mania et al., 2019, Wang and Janson, 2022]. There are also some connections between our work and the areas of model predictive control and system identification, but we defer these to the appendix (Appendix B) in the interest of space because the connections to our work are not as strong as the works surveyed in the rest of this subsection.

The two previous works that are most closely related to this paper are Dean et al. [2019] and Li et al. [2021], which both study safety-constrained LQR learning with unknown dynamics. Both works study the regret with respect to the baseline of the best linear controllers of the form $u_t = -Kx_t$ and derive an upper bound of $\tilde{O}_T(T^{2/3})$ on the regret. In Dean et al. [2019], they use system level synthesis to develop an algorithm that can safely inject noise into the system to give statistical guarantees on the learning rate. Li et al. [2021] provide the first adaptive learning algorithm for constrained LQR learning with unknown dynamics using a CE approach. While their results hold for higher dimensional LQR, our results improve on theirs in two ways. First, we are able to show a regret rate of $\tilde{O}_T(\sqrt{T})$ for some noise distributions, an improvement over their regret rate of $\tilde{O}_T(T^{2/3})$. Second, our regret results are with respect to a significantly stronger and more general baseline. These previous works focused on regret with respect to the best safe linear controller. However, the class of safe linear controllers is a relatively weak class of safe controllers, and the best safe linear controller can be far worse than the best overall safe controller. See Section 3.1 for more discussion on the importance of the choice of baseline. Note that these works allow constraints on both control and positions, while our results focus only on positional constraints. See Section 5 for more discussion on control constraints.

This work is the first part of a two part series of papers on safe LQR learning. In this paper, we provide a general framework for studying baselines of non-linear controllers in this problem. The second part of the series [Schiffer and Janson, 2025] uses the general framework from this paper to study a specific baseline of non-linear controllers known as the truncated linear controllers. Schiffer and Janson [2025] shows that the assumptions proposed in this paper do actually hold for a very natural baseline, and therefore our general framework can be applied to that baseline. Using the current paper's framework, Schiffer and Janson [2025] shows an even stronger result holds for that paper's specific baseline using a more complex algorithm. The additional results in Schiffer and Janson [2025] required significant extra technical work that would have pushed this paper's already lengthy technical appendix to an inaccessible length, which is why we split the results into a two part series.

# 2 Preliminaries

## 2.1 Outline of Preliminaries

In order to formally state our problem, the preliminaries section is organized as follows. First, in Section 2.2, we outline the dynamics of the system and the notation we will use for controllers. In Section 2.3, we define and motivate the expected-position safety constraints

we use to represent safety throughout the paper. In order for it to be possible to learn safely, we also need some initial information. In Section 2.4, we outline the exact assumptions we make on the initial uncertainty. Finally, in Section 2.5, we put everything from the previous sections together with a definition of regret to formally state our problem.

## 2.2 Problem Dynamics

Denote the state of the system at time $t$ for $t \in [T]$ as $x_t \in \mathbb{R}$ and the control at time $t$ as $u_t \in \mathbb{R}$. For simplicity, we will assume that the system starts at position $x_0 = 0$. The position at time $t + 1$ follows dynamics $x_{t+1} = a^* x_t + b^* u_t + w_t$, where $a^* \in \mathbb{R}$ and $b^* \in \mathbb{R}$ determine the dynamics and $w_t \overset{\text{i.i.d.}}{\sim} \mathcal{D}$ is the noise term drawn from a continuous, mean-0 probability distribution $\mathcal{D}$ with cumulative distribution function $F_{\mathcal{D}}$ and variance $\sigma_{\mathcal{D}}^2 = 1$. We will consider the quadratic cost at time $t$ as $q x_t^2 + r u_t^2$ for $q, r \in \mathbb{R}_{>0}$, and consider the sum of cost over the first $T$ steps. Throughout this paper, we will assume that the dynamics $a^*, b^*$ are unknown, while all other problem parameters are known (e.g. $\mathcal{D}, q, r$, etc.). For simplicity, we will denote the unknown dynamics as $\theta^* = (a^*, b^*) \in \mathbb{R}^2$.

We will also use the following controller notation. Define $H_t = (x_0, u_0, x_1, ..., u_{t-1}, x_t)$, and $\mathcal{F}_t = \sigma(H_t)$, the sigma algebra generated by $H_t$. We define a (possibly time-dependent and randomized) controller $C$ such that the control chosen at time $t$ is $u_t = C(H_t)$. Note that any randomness in the controller $C$ must be independent of the noise random variables $\{w_t\}_{t=0}^{T-1}$. Define the $T$-step *cost* of a controller $C$ starting at position $x_0$ under dynamics $\theta$ with noise random variables $W = \{w_t\}_{t=0}^{T-1}$ as

$$J(\theta, C, T, x_0, W) = \frac{1}{T} \left( q x_T^2 + \sum_{t=0}^{T-1} q x_t^2 + r u_t^2 \right), \tag{1}$$

$$\text{where } u_t = C(H_t), \ x_{t+1} = a x_t + b u_t + w_t, \ w_t \overset{\text{i.i.d.}}{\sim} \mathcal{D}.$$

Notice that $J$ outputs an average cost. We will refer to $T \cdot J(\theta, C, T, x_0, W)$ as the *total cost*. We denote $J^*(\theta, C, T, x_0)$ as the expectation of $J(\theta, C, T, x_0, W)$ with respect to only the randomness in $W$. Formally, this means that $J^*(\theta, C, T, x_0) = \mathbb{E}\left[J(\theta, C, T, x_0, W) \mid \theta, C, T, x_0\right]$ in case any of $\theta, C, T$, and $x_0$ are random, but in the typical setting when $\theta, C, T$, and $x_0$ are all deterministic, $J^*(\theta, C, T, x_0)$ will be non-random. For notational simplicity, we also define $J^*(\theta, C, T) = J^*(\theta, C, T, 0)$.

## 2.3 Constraints

Now we will formalize our positional constraints. Both Dean et al. [2019] and Li et al. [2021] formulate their positional constraints as *realized-position constraints* of the form

$$D_{\text{L}}^x \leq x_t \leq D_{\text{U}}^x, \tag{2}$$

which must be satisfied with probability 1 when the dynamics are known. Realized-position constraints that hold with probability 1 have the easy interpretation that the realized position must never exceed the realized-position boundaries given by the user of the algorithm.

However, in the case of unbounded noise distributions (for example Gaussian noise), having the realized position never exceed any compact set with probability 1 is impossible even with known dynamics. This is because with Gaussian noise, there is always a strictly positive probability that $x_t$ will be outside of the safe region $[D_{\mathrm{L}}^x, D_{\mathrm{U}}^x]$ for any choice of control $u_{t-1}$. Therefore, in order to allow for unbounded noise distributions, we must relax the requirement of never exceeding the constraints with probability 1, and instead allow the position $x_t$ to exceed the realized-position boundaries $D_{\mathrm{L}}^x$ and $D_{\mathrm{U}}^x$ with probability at most $\delta_{\mathrm{traj}}$ throughout the entire trajectory. Using a union bound, one way to achieve this relaxation for $T$ steps is to require that for every $t$,

$$D_{\mathrm{L}}^x - F_{\mathcal{D}}^{-1}\left(\frac{\delta_{\mathrm{traj}}}{2T}\right) \le a^* x_t + b^* u_t \le D_{\mathrm{U}}^x - F_{\mathcal{D}}^{-1}\left(1 - \frac{\delta_{\mathrm{traj}}}{2T}\right). \tag{3}$$

Motivated by this result, we will formulate our problem in terms of *expected-position constraints* of the form

$$D_{\mathrm{L}}^{\mathbb{E}[x]} \le a^* x_t + b^* u_t \le D_{\mathrm{U}}^{\mathbb{E}[x]}. \tag{4}$$

Because $\mathcal{D}$ is mean-0, this expected-position constraint has the easy interpretation of constraining the expected position, conditional on the history, at every time point (hence the $\mathbb{E}[x]$ superscript). By constraining the expected position, we are also implicitly constraining the realized position $x_t$ to be within the random interval $[D_{\mathrm{L}}^{\mathbb{E}[x]} + w_{t-1}, D_{\mathrm{U}}^{\mathbb{E}[x]} + w_{t-1}]$. Furthermore, if the noise distribution has support $[-\bar{w}, \bar{w}]$ and $\delta_{\mathrm{traj}} = 0$ (as in Dean et al. [2019] and Li et al. [2021]), then realized-position constraints are a special case of the expected-position constraints: Equation (2) with realized-position boundaries $D^x := \left(D_{\mathrm{L}}^x, D_{\mathrm{U}}^x\right)$ is equivalent to Equation (4) with expected-position boundaries $D^{\mathbb{E}[x]} := (D_{\mathrm{L}}^{\mathbb{E}[x]}, D_{\mathrm{U}}^{\mathbb{E}[x]}) = \left(D_{\mathrm{L}}^x + \bar{w}, D_{\mathrm{U}}^x - \bar{w}\right)$. For unbounded noise, Equation (2) is impossible to satisfy with probability 1, while Equation (4) is possible to satisfy and is directly related to the problem of satisfying the realized-position constraints with high probability. Therefore, the constraints in Equation (4) can in some sense be thought of as a generalization of the realized-position constraints in Equation (2). To maintain that 0 is a safe position, we will also require that $D_{\mathrm{L}}^{\mathbb{E}[x]} < 0 < D_{\mathrm{U}}^{\mathbb{E}[x]}$ (see Assumption 3).

In order to satisfy the realized position constraints in Equation (2) for all $T$ steps with constant probability, the magnitude of the boundaries must scale with the max position, which scales with the magnitude of the largest realized noise. For an unbounded distribution $\mathcal{D}$, this means that the realized-position boundaries must be a function of $T$ that grows with $T$. Now looking at Equation (3), the implied expected-position constraints include both $D^x$ (which may be a function of $T$) and a quantile of the noise distribution (which is explicitly a function of $T$). Therefore, we will allow the expected-position boundary $D^{\mathbb{E}[x]}$ of Equation (4) to depend on $T$. However, in the typical feasible safe RL problem we will have expected-position boundaries that are $O_T(1)$. The reason for this is that the expected-position constraints only bound the position in expectation. Therefore, unlike the realized-position boundary which must scale with the maximum position in order to be feasible, the expected-position boundary is feasible as long as it scales with the largest product of position and dynamics estimation error (uncertainty in $\theta$). Under the assumptions in this paper, we will achieve an estimation error that decreases at a rate that is much faster than the rate at which the maximum position grows. Thus, while we allow the expected-position boundaries

to be functions of $T$, the reader should generally think of them as not growing with $T$ in a typical problem, and indeed some of our results will explicitly require the expected-position boundaries to be $O_T(1)$.

Formally, we define safety as follows. Note that when the boundaries $(D_L^{\mathbb{E}[x]}, D_U^{\mathbb{E}[x]})$ are clear in context, we will drop the constraints and simply refer to algorithms that are safe for a specific dynamics $\theta^*$.

**Definition 1.** *A series of controls $\{u_t\}_{t=0}^{T-1}$ are safe for dynamics $\theta^*$ and boundaries $(D_L^{\mathbb{E}[x]}, D_U^{\mathbb{E}[x]})$ if every control satisfies Equation (4). Similarly, a controller $C$ is safe for dynamics $\theta^*$ and boundaries $(D_L^{\mathbb{E}[x]}, D_U^{\mathbb{E}[x]})$ if the resulting controls $\{C(H_t)\}_{t=0}^{T-1}$ under true dynamics $\theta^*$ are safe for dynamics $\theta^*$.*

## 2.4 Initial Uncertainty Assumptions

Without any prior knowledge about the unknown dynamics $\theta^*$, it is impossible to choose a first action that is guaranteed to be safe for all $\theta^* \in \mathbb{R}^2$. Therefore, to learn anything about the unknown dynamics while maintaining safety, we require some initial information about the unknown dynamics. Before getting into our main results, we will therefore formalize our assumptions about the initial uncertainty in our problem. As is standard in previous works [Abbasi-Yadkori and Szepesvári, 2011, Li et al., 2021], we will assume the following:

**Assumption 1.** *The algorithm has access to some $\Theta = \Theta_a \times \Theta_b = [\underline{a}, \bar{a}] \times [\underline{b}, \bar{b}]$ such that $\theta^* \in \Theta$ and $\bar{b} \geq \underline{b} > 0$ and $\bar{a} \geq \underline{a} > 0$.*

$\Theta$ can be thought of as the initial uncertainty set for $\theta^*$. Define the size of such a set $\Theta$ as $\text{size}(\Theta) = \max(\bar{a} - \underline{a}, \bar{b} - \underline{b})$. Note that depending on the size of $\Theta$, maintaining safety with respect to the expected-position boundaries for any $\theta^* \in \Theta$ may be infeasible. Infeasible in our setting means that there does not exist any adaptive controller $C$ such that for all $\theta^* \in \Theta$, the controller is safe with high probability, i.e. $\mathbb{P}\left(\forall t < T : D_L^{\mathbb{E}[x]} \leq a^* x_t + b^* C(H_t) \leq D_U^{\mathbb{E}[x]}\right) \geq 1 - \delta$. Clearly feasibility of $\Theta$ (for some appropriate choice of $\delta$) is a necessary condition for our problem to have a solution. The assumptions we make are only slightly stronger than just feasibility, which we discuss further in Appendix H.3. As described in Section 1.4, many previous works have developed algorithms that maintain guaranteed safety, but to the best of our knowledge the exact amount of prior information needed has not been quantified.

The assumption that $a^*, b^* > 0$ is for algebraic convenience, and the same results can be shown for any constant $a^*, b^* \in \mathbb{R}$. The assumption that $\underline{a}, \underline{b} > 0$ can actually be removed given the next assumption, and we discuss this more in Appendix H.1.

The other main assumption about prior information that we make is that we have sufficient information to not violate the safety constraint for some initial period of the algorithm.

**Assumption 2.** *There is a known controller $C^{\text{init}}$ such that $\forall x \in \left[D_L^{\mathbb{E}[x]} + F_{\mathcal{D}}^{-1}(\frac{1}{T^4}), D_U^{\mathbb{E}[x]} + F_{\mathcal{D}}^{-1}(1 - \frac{1}{T^4})\right]$,*

$$D_L^{\mathbb{E}[x]} + \frac{b^*}{\log(T)} \leq a^* x + b^* C^{\text{init}}(x) \leq D_U^{\mathbb{E}[x]} - \frac{b^*}{\log(T)}. \tag{5}$$

8

To get a sense of how strong Assumption 2 is, note that if we ignore the vanishing log terms in Equation (5), then Assumption 2 is equivalent to assuming that we can identify any safe controller. If this is not the case, then safe learning is clearly impossible. We further discuss Assumption 2 and how it relates to the concept of feasibility in Appendix H.3. In Appendix H.2, we also provide further interpretation of Assumption 2 in the case of bounded noise.

## 2.5   Problem Statement

We define $\mathcal{C}^{\theta^*}$ as a baseline class of controllers if every controller $C \in \mathcal{C}^{\theta^*}$ is safe with respect to dynamics $\theta^*$ with probability 1. **If $\theta^*$ were known**, then the safe LQR problem with $\mathcal{C}^{\theta^*}$ as the baseline would simply be to minimize the expected total cost for all controllers in this baseline, i.e. to solve

$$\min_{C \in \mathcal{C}^{\theta^*}} \quad T \cdot J^*(\theta^*, C, T). \tag{6}$$

We will use the expression in Equation (6) as the baseline cost to which we compare the cost of our algorithms. We will often consider families of controller classes $\{\mathcal{C}^\theta\}_{\theta \in \Theta}$ such that for any dynamics $\theta$, every controller in the class $\mathcal{C}^\theta$ is safe for dynamics $\theta$ with probability 1. For example, the baseline class $\mathcal{C}^\theta$ could be the class of linear controllers that are safe for dynamics $\theta$, the class of affine controllers that are safe for dynamics $\theta$, all controllers that are safe for dynamics $\theta$, etc.

The *regret* of an algorithm with corresponding controller $C_{\mathrm{alg}}$ with respect to baseline $\mathcal{C}^{\theta^*}$ is the random variable

$$\text{Regret} \ := T \cdot J(\theta, C_{\mathrm{alg}}, T, 0, W) - \min_{C \in \mathcal{C}^{\theta^*}} T \cdot J^*(\theta^*, C, T). \tag{7}$$

Note that this regret random variable compares the realized cost of the algorithm with the expected cost of a controller from the baseline class, and this definition of regret is typical in the LQR learning literature [Abbasi-Yadkori and Szepesvári, 2011, Li et al., 2021]. We also could have defined regret comparing the realized cost of an algorithm to the realized cost of the best (in expectation) controller from the baseline class. Due to standard concentration inequalities, the realized total cost of the baseline controller will be within $\tilde{O}(\sqrt{T})$ of the expected total cost of the baseline controller. Therefore, considering a realized total cost for both terms in the regret would change our regret bounds by at most $\tilde{O}(\sqrt{T})$ and therefore not change any of the results.

The overarching goal of this paper is to find a controller $C_{\mathrm{alg}}$ that achieves low regret as defined in Equation (7) and such that for any true dynamics $\theta^* \in \Theta$, the controller $C^{\mathrm{alg}}$ is safe for $\theta^*$ with probability $1 - o_T(1/T)$. Note that we only require that the algorithm $C_{\mathrm{alg}}$ is safe with probability $1 - o_T(1/T)$, while we require the baseline to be safe with probability 1. This (slightly unfair) mismatch is necessary to allow the algorithm to use information "learned" from historical observations when trying to satisfy the safety constraints. For example, if $\mathcal{D}$ is an unbounded distribution, then it is impossible to conclude anything with probability 1 based on any amount of historical information. We want to allow our algorithm to use information about $\theta^*$ learned from previous time steps to choose better future safe controls, and therefore we only require safety with respect to $\theta^*$ with probability $1 - o_T(1/T)$. We

chose $1 - o_T(1/T)$ for the safety probability because this is strictly stronger than $1 - o_T(1)$ or $1 - \delta$ for constant $\delta > 0$, and therefore our results hold for these larger probabilities of satisfying safety as well. In principle, we could also compare to a baseline that allows some probability of error. However, because the baseline does not need to learn $\theta^*$, allowing it to be safe with probability slightly less than 1 would not significantly impact its cost, while it would significantly increase the mathematical complexity of the analysis.

Finally, we will make the following assumptions about the problem specifications throughout this paper.

**Assumption 3** (Problem Specifications). *The noise distribution $\mathcal{D}$ is mean-0, variance 1, and subgaussian with bounded density. The boundaries $D_{\mathrm{L}}^{\mathbb{E}[x]}, D_{\mathrm{U}}^{\mathbb{E}[x]}$ (which may be functions of $T$) satisfy that $-\log^2(T) \leq D_L^{\mathbb{E}[x]} < 0 < D_U^{\mathbb{E}[x]} \leq \log^2(T)$ and that $D_U^{\mathbb{E}[x]} - D_L^{\mathbb{E}[x]} \geq \frac{1}{\log(T)}$.*

For exposition purposes, we also assume that $\log_2(T^{1/12})$ is an integer. The assumption of variance 1 gives a simpler uncertainty bound, but as in Abbasi-Yadkori and Szepesvári [2011] this can be relaxed. We assume that $\max(|D_{\mathrm{L}}^{\mathbb{E}[x]}|, D_{\mathrm{U}}^{\mathbb{E}[x]}) \leq \log^2(T)$ because if the constraints are greater than $\log^2(T)$, then the constraints have very little impact on the optimal controller. This is because with subgaussian noise, with high probability the noise random variables have magnitude less than $o(\log(T))$, and so reasonable controllers will with high probability never hit the constraint. Therefore, if both boundaries are greater than $\log^2(T)$ then the problem becomes similar to the unconstrained problem, and if one boundary is large, then the problem becomes one sided which is an easier version of our problem. The assumption of mean-0 and subgaussian noise is also standard in the LQR literature [Abbasi-Yadkori and Szepesvári, 2011, Dean et al., 2019, Li et al., 2021].

Putting everything together, the formal problem we are considering is the following.

**Problem 1** (Safe LQR Learning). *Suppose we are given $D, \mathcal{D}, \Theta, T$ that satisfy Assumption 1–3 and a set of baseline classes of controllers $\{\mathcal{C}^\theta\}_{\theta \in \Theta}$. Then the goal of safe LQR learning is to find an algorithm $C^{\mathrm{alg}}$ that achieves a regret with respect to baseline $\mathcal{C}^{\theta^*}$ that is as low as possible, while also satisfying $\sup_{\theta \in \Theta} \mathbb{P}\left(C^{\mathrm{alg}} \text{ is safe with respect to } \theta\right) = 1 - o_T(1/T)$.*

Note that $\sup_{\theta \in \Theta} \mathbb{P}\left(C^{\mathrm{alg}} \text{ is safe with respect to } \theta\right) = 1 - o_T(1/T)$ is equivalent to requiring that there exists some probability $p = 1 - o_T(1/T)$ such that for any $\theta \in \Theta$, if the true dynamics $\theta^* = \theta$ then the controls used by $C^{\mathrm{alg}}$ are safe with respect $\theta^*$ with probability $p$.

## 2.6 Notation

To simplify notation, we use $\theta = (a, b)$ to represent an arbitrary set of dynamics and $\theta^* = (a^*, b^*)$ to represent the true (unknown) dynamics. We will also use $D := (D_{\mathrm{L}}, D_{\mathrm{U}}) := (D_{\mathrm{L}}^{\mathbb{E}[x]}, D_{\mathrm{U}}^{\mathbb{E}[x]})$ (i.e., drop the superscripts). We will use $\tilde{O}_T$ and $O_T$ notation to represent $\tilde{O}$ and $O$ with respect to $T$, where the values of the hidden constants and log terms may depend on the values of problem inputs such as $q, r, D, \mathcal{D}, \Theta$. Because the nature of our problem requires us to define a significant amount of notation in this paper, we have a table in Appendix A that lists the common notation used throughout the paper that the reader can use as a reference if needed.

# 3    Theoretical Results

The goal of this paper is to provide a general framework for studying the regret with respect to non-linear baselines of controllers. We first introduce a general class of baselines satisfying regularity conditions in Section 3.1. We then present our two main theorems in Section 3.2.

## 3.1    Regret Rates for General Baselines

In order to present our main theorem, we first need a baseline class of controllers $\mathcal{C}^{\theta^*}$ to define the regret in Equation (7). In both Li et al. [2021] and Dean et al. [2019], the regret baseline for the $\tilde{O}_T(T^{2/3})$ results is the cost of the best stationary linear controller of the form $u_t = -Kx_t$ that is safe for $\theta^*$ with probability 1. We will refer to the class of stationary linear controllers that are safe for $\theta^*$ with probability 1 as the class of safe linear controllers. Since not all linear controllers are safe for dynamics $\theta^*$, this is restricted to $K$ that will maintain safety for $\theta^*$ for any realization of the noise, and therefore can be a very weak baseline. Linear controllers are not always well-suited for safety constrained LQR because linear controllers only have one degree of freedom $(K)$, but in safety constrained LQR the controller must balance keeping regret low with being safe. For example, when $D_{\mathrm{U}}$ and $D_{\mathrm{L}}$ are not symmetric, the best linear controller must still behave symmetrically. However, symmetric behavior may be far from optimal for $D_{\mathrm{U}}$ and $D_{\mathrm{L}}$ that are not symmetric, and linear controllers lack the flexibility to behave non-symmetrically. Therefore, there exist much stronger baselines than the safe linear controllers studied in Li et al. [2021], Dean et al. [2019].

In Section 3.2, we present two results that hold for a wide range of stronger baseline classes of controllers. Before stating the theorems, we will outline a few assumptions on the controllers in these general baseline classes.

Let $\{\mathcal{C}^{\theta}\}_{\theta \in \Theta}$ be the set of baseline classes of controllers for dynamics $\theta \in \Theta$. For the rest of this paper, we will assume that the baseline class of controllers satisfies Assumption 4.

**Assumption 4.** *All of the controllers in the baseline class $\mathcal{C}^{\theta}$ for all $\theta \in \Theta$ are stationary, Markovian, deterministic, and safe for dynamics $\theta$ with probability 1.*

Note that the assumption that every controller in $\mathcal{C}^{\theta}$ is safe for dynamics $\theta$ with probability 1 is consistent with the baselines of Li et al. [2021] and Dean et al. [2019]. Additionally, this means that the baseline class of controllers could change depending on the dynamics $\theta$, as the class of controllers that is safe for one dynamics will not necessarily be safe for a different dynamics. One option is to construct the baseline class from another class of controllers $\tilde{\mathcal{C}}$ (for example the class of linear controllers), as follows:

$$\{C \in \tilde{\mathcal{C}} : C \text{ is safe for dynamics } \theta\}. \tag{8}$$

If $\tilde{\mathcal{C}}$ is a rich enough class of controllers (e.g. all controllers), then Equation (8) would result in a good safe baseline. However, if $\tilde{\mathcal{C}}$ is a relatively small class of controllers (e.g. linear controllers), then the restriction in Equation (8) to only controllers in the class that are safe for $\theta$ may result in a weak safe baseline. Therefore, instead of simply subsetting the class of controllers $\tilde{\mathcal{C}}$ as in Equation (8), we will preserve the complexity of the function class $\tilde{\mathcal{C}}$

by transforming *every* controller in $\tilde{\mathcal{C}}$ into a controller that is safe for $\theta$. We generalize even further by allowing the starting class of controllers $\tilde{\mathcal{C}}^\theta$ to be different for each $\theta$.

**Assumption 5** (Truncation). *For any $\theta$, there exists a controller class $\tilde{\mathcal{C}}^\theta$ of deterministic controllers such that the baseline class $\mathcal{C}^\theta$ consists of all controllers of the following form for $C \in \tilde{\mathcal{C}}^\theta$:*

$$
C^\theta(x) = \begin{cases} C(x) & \text{if } D_{\mathrm{L}} \leq ax + bC(x) \leq D_{\mathrm{U}} \\ \frac{D_{\mathrm{U}} - ax}{b} & \text{if } ax + bC(x) > D_{\mathrm{U}} \\ \frac{D_{\mathrm{L}} - ax}{b} & \text{if } ax + bC(x) < D_{\mathrm{L}}. \end{cases} \tag{9}
$$

By this construction, every controller $C^\theta \in \mathcal{C}^\theta$ is safe for dynamics $\theta$. We will also assume that $\mathcal{C}^\theta$ is parameterizable by a scalar parameter $K \in \mathbb{R}$. This allows us to choose the optimal controller in $\mathcal{C}^\theta$ in terms of the parameter $K$.

**Assumption 6** (Parametrization). *For any $\theta$, there exists $K_{\mathrm{L}}^\theta, K_{\mathrm{U}}^\theta \in \mathbb{R}$ such that the $\mathcal{C}^\theta$ in Assumption 5 can be parameterized as $\mathcal{C}^\theta = \{C_K^\theta : K \in [K_{\mathrm{L}}^\theta, K_{\mathrm{U}}^\theta]\}$. Furthermore, for any $\theta$, $T$ there exists a $K_{\mathrm{opt}}(\theta, T)$ satisfying*

$$
K_{\mathrm{opt}}(\theta, T) = \arg \min_{K \in [K_{\mathrm{L}}^\theta, K_{\mathrm{U}}^\theta]} J^*(\theta, C_K^\theta, T).
$$

Our results require two more key assumptions on the class of controllers.

**Assumption 7** (Average Cost Lipschitz in Optimal Controller). *There exists $\epsilon_{\mathrm{A7}} = \tilde{\Omega}_T(1)$ such that for any $\|\theta - \theta^*\|_\infty \leq \epsilon_{\mathrm{A7}}$ and $t \leq T$,*

$$
|J^*(\theta^*, C_{K_{\mathrm{opt}}(\theta,t)}^\theta, t) - J^*(\theta^*, C_{K_{\mathrm{opt}}(\theta^*,t)}^{\theta^*}, t)| \leq \tilde{O}_T\left(\|\theta - \theta^*\|_\infty + \frac{1}{T^2}\right).
$$

Assumption 7 relates the expected cost under dynamics $\theta^*$ of the optimal controller for dynamics $\theta^*$ to the expected cost of the optimal controller for some other dynamics $\theta$ close to $\theta^*$. Intuitively, this is a form of Lipschitz continuity which implies that the performance of the optimal controller is not too sensitive to the choice of $\theta$. Some sort of continuity assumption is intuitively required for any form of certainty equivalence algorithm to achieve low regret guarantees.

**Assumption 8** (Total Cost Lipschitz in Initial Position). *There exist $\epsilon_{\mathrm{A8}}, \delta_{\mathrm{A8}} = \tilde{\Omega}_T(1)$ such that for any $\theta$ satisfying $\|\theta - \theta^*\|_\infty \leq \epsilon_{\mathrm{A8}}$ the following holds. For $t < T$, let $W' = \{w_i\}_{i=0}^{t-1}$. Then for any $K \in [K_{\mathrm{L}}^\theta, K_{\mathrm{U}}^\theta]$, there exists a set $\mathcal{Y}_{\mathrm{A8}} \in \mathbb{R}^t$ that depends only on $C_K^\theta$ such that the following holds. Define $E_{\mathrm{A8}}\left(C_K^\theta, W'\right)$ as the event that $W' \in \mathcal{Y}_{\mathrm{A8}}$. Then $\mathbb{P}(E_{\mathrm{A8}}\left(C_K^\theta, W'\right)) \geq 1 - o_T(1/T^{10})$ and for any $|x|, |y| \leq 4\log^2(T)$ such that $|x - y| \leq \delta_{\mathrm{A8}}$, conditional on event $E_{\mathrm{A8}}\left(C_K^\theta, W'\right)$,*

$$
\left|t \cdot J(\theta^*, C_K^\theta, t, x, W') - t \cdot J(\theta^*, C_K^\theta, t, y, W')\right| \leq \tilde{O}_T(|x - y| + \|\theta - \theta^*\|_\infty). \tag{10}
$$

Assumption 8 relates the random variables of cost when starting at two different positions, $x$ and $y$, but with the same noise random variables $W'$. Intuitively, this implies that making a small non-optimal control will not have significant long-term impact on the total cost.

Therefore, this assumption can be thought of as assuming the total cost is Lipschitz in the initial position.

The final assumption we consider is an assumption that the noise distribution has sufficiently large support, which we require for Theorem 1 but not for Theorem 2. Note that we only need the noise distribution to have large support relative to one of the two boundaries ($D_\mathrm{U}$ or $D_\mathrm{L}$). We will w.l.o.g. state Assumption 9 relative to boundary $D_\mathrm{U}$, however an equivalent assumption swapping $D_\mathrm{L}$ and $D_\mathrm{U}$ would also be sufficient for Theorem 1.

**Assumption 9.** *For any $z$, define $P(\theta, K, z)$ as the largest $x$ such that $ax + bC_K^\theta(x) \leq z$. There exists a constant $\epsilon_\mathrm{A9} > 0$ such that the following equation holds for all $t \geq \sqrt{T}$ for sufficiently large $T$:*

$$\mathbb{P}_{w \sim \mathcal{D}} \left( w \geq P(\theta^*, K_\mathrm{opt}(\theta^*, t), D_\mathrm{U}) - D_\mathrm{L} \right) \geq \epsilon_\mathrm{A9} > 0. \tag{11}$$

The quantity $P(\theta^*, K_\mathrm{opt}(\theta^*, t), D_\mathrm{U})$ will often be proportional to and greater than $D_\mathrm{U}$. Because $\mathcal{D}$ is constant relative to $T$, Assumption 9 implies that the boundary $D$ must satisfy $\|D\|_\infty = O_T(1)$. When $\|D\|_\infty = O_T(1)$, Assumption 9 effectively requires that the noise distribution $\mathcal{D}$ has a constant probability of spanning the distance between $D_\mathrm{L}$ and $D_\mathrm{U}$. Note that Assumption 9 is automatically satisfied for any $\|D\|_\infty = O_T(1)$ when the noise distribution is Gaussian, unbounded, or bounded with a high enough variance. This assumption will be necessary to achieve regret of $\tilde{O}_T(\sqrt{T})$ in Theorem 1, as the variance from the noise distribution of Assumption 9 provides the controller with extra exploration that leads to better estimation. We will also provide a result for general classes of controllers that does not require this assumption, but achieves a worse regret rate (Theorem 2).

## 3.2 Theorems

We are now ready to present our first theorem.

**Theorem 1.** *In the setting of Problem 1 and under further Assumptions 4–9, there exists an algorithm $C^\mathrm{alg}$ (Algorithm 3) that with probability $1 - o_T(1/T)$ achieves $\tilde{O}_T(\sqrt{T})$ regret with respect to baseline $\mathcal{C}^{\theta^*}$ while also satisfying $\sup_{\theta \in \Theta} \mathbb{P} \left( C^\mathrm{alg} \text{ is safe with respect to } \theta \right) = 1 - o_T(1/T)$.*

The key lemma in proving Theorem 1 is a new estimation bound for the unknown system dynamics $\theta^*$ (Lemma 26). Informally, this estimation bound shows that simply by obeying safety constraints, the unknown dynamics can be estimated at a rate of $1/\sqrt{t}$ without injecting any additional randomness into the controller. This faster rate of learning is because in order to be safe, the controller must frequently be non-linear, which in turn helps learn the unknown dynamics. This result of safe behavior leading to faster learning rates may also be of independent interest in other safe RL problems.

The more general result of this paper is Theorem 2, which achieves a weaker regret rate of $\tilde{O}_T(T^{2/3})$ but applies for any subgaussian noise distribution (in particular, it drops Assumption 9).

**Theorem 2.** *In the setting of Problem 1 and under further Assumptions 4–8, there exists an algorithm $C^\mathrm{alg}$ (Algorithm 2) that with probability $1 - o_T(1/T)$ achieves $\tilde{O}_T(T^{2/3})$ regret*

with respect to baseline $\mathcal{C}^{\theta^*}$ while also satisfying $\sup_{\theta \in \Theta} \mathbb{P} \left( C^{\mathrm{alg}} \text{ is safe with respect to } \theta \right) = 1 - o_T(1/T)$.

Theorem 2 is an improvement on existing results in that it bounds the regret of constrained LQR learning for any subgaussian noise distribution. See Section 4.1 and Section 4.2 for the proof sketches of Theorem 2 and Theorem 1 respectively. Previous works focus on linear controller baselines, and linear controllers have properties that allow for easier regret analysis. Theorems 2 and 1 reduce these "useful" properties of linear controllers to Assumptions 7 and 8. Therefore, many classes of non-linear controllers can be constructed as described in this section, and all that needs to be done to show that the result of the theorems hold with such a class of controllers as a baseline is to show that this class of controllers satisfies Assumptions 7 and 8. Both of Assumptions 7 and 8 are simply Lipschitz conditions on the cost function (one with respect to the optimal controller and one with respect to the starting position), and therefore are likely to hold for many classes of controllers. In particular, Schiffer and Janson [2025] shows that both of these assumptions are satisfied for the class of truncated linear controllers, and therefore Theorems 1 and 2 apply for this baseline class of controllers. The properties in Assumptions 7 and 8 are the main tools that allow us to analyze the regret of nonlinear general baselines, and therefore these properties may be of independent interest outside of these theorems.

The algorithms that achieve the regret bounds of Theorems 1 and 2 follow the same general form. We outline the algorithm that achieves Theorem 2 below in Algorithm 1.

---

**Algorithm 1** Outline of Algorithm 2 for proof of Theorem 2

---

1: Explore for $\tilde{\Theta}_T(T^{2/3})$ steps using controller $C^{\mathrm{init}}$ from Assumption 2 with random noise.
2: **for** $s \in [0 : \log(T^{1/3}) - 1]$ **do**
3:      $\hat{\theta}_s \leftarrow$ regularized least-squares estimate of $\theta^*$ using data seen so far
4:      $\epsilon_s \leftarrow$ high probability bound on $\|\theta^* - \hat{\theta}_s\|_\infty$
5:      $C_s^{\mathrm{alg}} \leftarrow$ optimal controller from baseline class for dynamics $\hat{\theta}_s$
6:      For next $T^{2/3}2^s$ steps, use controller $C_s^{\mathrm{alg}}$ modified at each step to be safe for all dynamics $\theta$ satisfying $\|\theta - \hat{\theta}_s\|_\infty \le \epsilon_s$

---

This algorithm mostly behaves like a standard certainty equivalence algorithm, first calculating the regularized least-squares estimate of $\theta^*$ and then finding the best controller for this estimated dynamics. This algorithm deviates from standard certainty equivalence in the final line, where the algorithm enforces safety by modifying the controller $C_s^{\mathrm{alg}}$. Because $\theta^*$ with high probability satisfies $\|\theta^* - \hat{\theta}_s\|_\infty \le \epsilon_s$, the modification in the final line guarantees safety for dynamics $\theta^*$ with high probability. The bulk of the theoretical work in proving Theorem 2 is upper bounding the regret contributed by these safety modifications. Theorem 1 follows a similar pattern with a slightly more complicated choice of $\hat{\theta}_s$. In the setting of Theorem 1, the large support of the noise distribution leads to the controls used by controller $C_s^{\mathrm{alg}}$ being non-linear by a constant amount for a constant fraction of the steps. This non-linearity allows the algorithm to learn at a faster rate than in Theorem 2 and results in the lower regret bound of $\tilde{O}_T(\sqrt{T})$. Note also that the length of the exploration period and the number of steps in each round of the loop are chosen differently for Algorithm 3 than for Algorithm 2. See the proof sketches in the following section for more details.

# 4   Proof Sketches of Main Results

We will present the proof sketches (and formal proofs) of the main results in *reverse* of the order in which they were stated in the previous section. We present the proofs in this manner because the result of Theorem 2 is a weaker result in a more general setting. We therefore build off of this proof in the subsequent proof of Theorem 1 by strengthening the result of Theorem 2 in less general settings.

## 4.1   Proof Sketch of Theorem 2

The full proof of Theorem 2 can be found in Appendix C.

First we state Algorithm 2, which is the algorithm that achieves the guarantee of Theorem 2. But before presenting the algorithm, we need some additional notation. Fix a constant $\lambda > 0$. Define $z_t = (x_t, u_t)^\top$ and $V_t = \lambda I + \sum_{i=0}^{t-1} z_i z_i^\top$, where $I$ is the identity matrix. Define $X_t$ as the column vector $(x_1, ..., x_t)^\top$ and $Z_t$ as the matrix with rows $z_0^\top, ..., z_{t-1}^\top$. Define $B_t = \alpha\sqrt{\log(\det(V_t)) + \log(\lambda^2) + 2\log(T^2)} + \sqrt{\lambda}(\bar{a}^2 + \bar{b}^2) = \tilde{O}_T(1)$ where $\alpha$ is the subgaussian parameter of $\mathcal{D}$. The algorithm that achieves the regret bound of Theorem 2 is given as Algorithm 2.

**Algorithm 2 Intuition**   Algorithm 2 can be broken into two phases: a warm-up exploration phase (Lines 2–4) and a safe certainty equivalence phase (Lines 5–13). In the warm-up phase, the controls are random which allows for sufficient exploration and learning of the unknown dynamics. In the certainty equivalence phase, $\hat{\theta}_s$ is the regularized least-square estimate of $\theta^*$ based on the data seen so far. $\epsilon_s$ is an upper bound on the distance between $\hat{\theta}_s$ and $\theta^*$ that holds with high probability. $C_s^{\mathrm{alg}}$ is the optimal controller from the baseline class for dynamics $\hat{\theta}_s$. Because $C_s^{\mathrm{alg}}$ is not guaranteed to be safe for dynamics $\theta^*$, we calculate $u_t^{\mathrm{safeU}}$ and $u_t^{\mathrm{safeL}}$ which are respectively the largest and smallest possible controls that satisfy the constraints for all dynamics $\theta$ within $\epsilon_s$ distance of $\hat{\theta}_s$ (which will with high probability include $\theta^*$). We then censor the control $C_s^{\mathrm{alg}}(x_t)$ with these two controls to guarantee with high probability that the final chosen control is safe with respect to dynamics $\theta^*$. In order to show Theorem 2, we must show that with probability $1 - o_T(1/T)$, Algorithm 2 is safe with respect to $\theta^*$ and that Algorithm 2 has $\tilde{O}_T(T^{2/3})$ regret. To show the latter, we will decompose the regret into four main components and consider each separately.

**Algorithm 2** Safe LQR for General Baselines

**Input:** $D, \mathcal{D}, \Theta, C^{\text{init}}, \{\mathcal{C}^\theta\}_{\theta \in \Theta}, T, \lambda$

1: $\nu_T \leftarrow T^{-1/3}$
2: **for** $t \leftarrow 0$ to $\frac{1}{\nu_T^2} - 1$ **do**  $\qquad\qquad\qquad\qquad\qquad$ ▷ Safe warm-up exploration phase
3: $\qquad \phi_t \sim \text{Rademacher}(0.5)$
4: $\qquad$ Use control $u_t = C^{\text{init}}(x_t) + \frac{\phi_t}{\log(T)}$

5: **for** $s \leftarrow 0$ to $\log_2(T\nu_T^2) - 1$ **do**  $\qquad\qquad\qquad\qquad$ ▷ Safe certainty equivalence phase
6: $\qquad T_s \leftarrow \frac{2^s}{\nu_T^2}$
7: $\qquad \hat{\theta}_s \leftarrow (Z_{T_s}^\top Z_{T_s} + \lambda I)^{-1} Z_{T_s}^\top X_{T_s}$
8: $\qquad C_s^{\text{alg}} \leftarrow C_{K_{\text{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}$
9: $\qquad \epsilon_s \leftarrow B_{T_s} \sqrt{\frac{\max(V_{T_s}^{22}, V_{T_s}^{11})}{V_{T_s}^{11} V_{T_s}^{22} - (V_{T_s}^{12})^2}}$
10: $\qquad$ **for** $t \leftarrow T_s$ to $2T_s - 1$ **do**
11: $\qquad\qquad u_t^{\text{safeU}} \leftarrow \max \left\{ u : \max_{\|\theta - \hat{\theta}_s\|_\infty \leq \epsilon_s} ax_t + bu \leq D_{\text{U}} \right\}$
12: $\qquad\qquad u_t^{\text{safeL}} \leftarrow \min \left\{ u : \min_{\|\theta - \hat{\theta}_s\|_\infty \leq \epsilon_s} ax_t + bu \geq D_{\text{L}} \right\}$
13: $\qquad\qquad$ Use control $u_t = \max \left( \min \left( C_s^{\text{alg}}(x_t), u_t^{\text{safeU}} \right), u_t^{\text{safeL}} \right)$

**Safety of Algorithm 2** We begin with analyzing the safety of Algorithm 2. The first loop (warm-up exploration) of Algorithm 2 is safe with respect to dynamics $\theta^*$ as a result of Assumption 2. In the second loop (safe certainty equivalence), the control in Line 13 is chosen to enforce safety relative to all $\theta$ satisfying $\|\theta - \hat{\theta}_s\|_\infty \leq \epsilon_s$. By the choice of $\epsilon_s$, the true dynamics $\theta^*$ satisfy $\|\theta^* - \hat{\theta}_s\|_\infty \leq \epsilon_s$ for all $s$ with probability $1 - o_T(1/T)$ (Lemma 23). Therefore, the control applied in Line 13 is safe with respect to $\theta^*$ for all $t$ with probability $1 - o_T(1/T)$. Therefore, Algorithm 2 is safe with respect to $\theta^*$ with probability $1 - o_T(1/T)$.

**Regret from warm-up period** The first component of regret $(R_0)$ is the cost of the warm-up exploration phase, which is the first $1/\nu_T^2$ steps of the algorithm. Using Assumption 3, we can show that the positions and controls during this phase are with high probability bounded by $\tilde{O}_T(1)$ (Lemma 4). Therefore, the cost during this phase can be bounded by $\tilde{O}_T(1/\nu_T^2)$ (Proposition 3). Importantly, after this initial exploration phase, $\epsilon_s = \tilde{O}_T(\nu_T)$ with probability $1 - o_T(1/T)$ (Lemma 2). This is a result of the Rademacher random variables in the warm-up phase.

**Regret from certainty equivalence** The second source of regret $(R_1)$ comes from the certainty equivalence aspect of the algorithm. In other words, $R_1$ is the regret from the fact that $K_{\text{opt}}(\hat{\theta}_s, T_s)$ is the optimal controller for dynamics $\hat{\theta}_s$ and not for dynamics $\theta^*$. By Lemma 2 and Lemma 23, with high probability $\|\hat{\theta}_s - \theta^*\|_\infty \leq \epsilon_s = \tilde{O}_T(\nu_T)$, so by Assumption 7 the expected cost of using controller $C_{K_{\text{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}$ for $T_s$ steps is at most $\tilde{O}_T(T_s \|\hat{\theta}_s - \theta^*\|_\infty + 1/T)$ more than the expected cost of using $C_{K_{\text{opt}}(\theta^*, T_s)}^{\theta^*}$ for $T_s$ steps. Using the aforementioned bound

comparing $\hat{\theta}_s$ and $\theta^*$, this source of regret can therefore be upper-bounded by $\tilde{O}_T(T\nu_T)$ with probability $1 - o_T(1/T)$ (Proposition 4).

**Regret from deviation from expectation**   The third source of regret ($R_2$) comes from the fact that we defined regret as the difference between the cost of the algorithm (which is a random variable) and the expected cost of the best controller in the baseline class (which is nonrandom). To bound this regret term, we show that the cost of the algorithm concentrates within $\tilde{O}_T(\sqrt{T})$ of its expectation with probability $1 - o_T(1/T)$ (Proposition 5). For this result, we use a variant of McDiarmid's Inequality that applies to high probability events combined with Assumption 8 (Lemma 6).

**Regret from enforcing safety**   The final source of regret ($R_3$) is a result of the times the algorithm "enforces safety" on the controls by sometimes using controls $u_t^{\text{safeU}}$ and $u_t^{\text{safeL}}$. With probability $1 - o_T(1/T)$, when the algorithm enforces safety, the chosen $u_t$ differs from $C_s^{\text{alg}}(x_t)$ by $\tilde{O}_T(\epsilon_s)$ (Lemma 9). By Assumption 8 and Lemma 2, the small differences between $C_s^{\text{alg}}(x_t)$ and $u_t$ each increase the cost by at most $\tilde{O}_T(\nu_T)$ with probability $1 - o_T(1/T)$. Therefore, the total cost of enforcing safety with these controls is $\tilde{O}_T(\nu_T T)$ with probability $1 - o_T(1/T)$ (Proposition 6).

**Combining Regret Terms**   Putting these four sources of regret together, the total regret can be upper bounded as follows with probability $1 - o_T(1/T)$:

$$T \cdot J(\theta^*, C_{\text{alg}}, T, 0, W) - T \cdot J^*(\theta^*, C_{K_{\text{opt}}(\theta^*, T)}^{\theta^*}, T) \leq R_0 + R_1 + R_2 + R_3 = \tilde{O}_T\left(\sqrt{T} + T\nu_T + \frac{1}{\nu_T^2}\right) = \tilde{O}_T(T^{2/3}),$$
(12)

where the last line comes from the fact that $\nu_T = T^{-1/3}$. See Appendix C and Equation (31) for a formal description of these four sources of regret.

## 4.2   Proof Sketch of Theorem 1

The full proof of Theorem 1 can be found in Appendix F.

**Algorithm and Intuition**   The algorithm that achieves the regret result of Theorem 1 is Algorithm 3, which is very similar to Algorithm 2. Rather than restating the entire algorithm here, we defer the full algorithm to the appendix and instead highlight the main differences between Algorithm 3 and Algorithm 2. The first modification is that for Algorithm 3 we choose $\nu_T = T^{-1/4}$, which affects the lengths of the exploration and certainty equivalence periods. The second major difference is that we change how $\hat{\theta}_s$ is defined. Recall that in Algorithm 2, $\hat{\theta}_s$ is the regularized least-squares estimate of $\theta^*$. For this algorithm we instead denote the regularized least-squares estimate as

$$\hat{\theta}_s^{\text{pre}} = (Z_{T_s}^\top Z_{T_s} + \lambda I)^{-1} Z_{T_s}^\top X_{T_s}.$$
(13)

Recall the function $P$ defined in Assumption 9. We choose $\hat{\theta}_s$ as

$$\hat{\theta}_s = \underset{\|\hat{\theta}_s - \hat{\theta}_s^{\text{pre}}\|_\infty \leq \epsilon_s}{\arg\min} \quad \underset{\|\theta - \hat{\theta}_s^{\text{pre}}\|_\infty \leq \epsilon_s}{\min} P(\theta, K_{\text{opt}}(\hat{\theta}_s, T_s), D_U).$$
(14)

The choice of $\hat{\theta}_s$ described above is a technical way of ensuring that $C_s^{\text{alg}}$ does sufficient exploration, which in turn guarantees a faster learning rate of the unknown dynamics. The key difference between the proof of Theorem 1 and the proof of Theorem 2 is a new upper bound on $\epsilon_s$ which is stronger than Lemma 2. Instead of $\epsilon_s = \tilde{O}_T(\nu_T)$ with probability $1 - o_T(1/T)$, we show that $\epsilon_s = \tilde{O}_T\left(\frac{1}{\sqrt{T_s}}\right)$ with probability $1 - o_T(1/T)$ (Lemma 19). Informally, this means that with high probability, the estimated dynamics at time $t$ are at most $\tilde{O}_T\left(\frac{1}{\sqrt{t}}\right)$ different from $\theta^*$, and this is a faster learning rate than in Theorem 2. This faster learning rate gives better upper bounds on the regret terms than in Theorem 2.

**Faster Learning Rate**  Showing the faster learning rate requires two main results. The first result is that the uncertainty $\epsilon_s$ can be upper-bounded by $\tilde{O}_T(1/\sqrt{|S_{T_s}|})$, where $|S_{T_s}|$ is the number of times Algorithm 3 uses the control $u_t^{\text{safeU}}$ before time $T_s$ (Lemma 21). To prove this result, we prove a more general uncertainty bound in Lemma 26. The key insight is that in order to maintain safety, the control $u_t^{\text{safeU}}$ will with high probability be sufficiently non-linear. This non-linearity combined with the variance in the position leads to a faster convergence rate of the upper bound in Lemma 23. The second result is that Algorithm 3 uses the control $u_t^{\text{safeU}}$ at least $\Omega_T(T_s)$ times before time $T_s$ (Lemma 20). The key insight to this result is that every time the position exceeds $P(\theta^*, K_{\text{opt}}(\theta^*, T_s), D_{\text{U}})$, Algorithm 3 will use control $u_t^{\text{safeU}}$. Assumption 9 says that the noise is large enough that (due to the choice of $\hat{\theta}_s$ in Equation (14)) the position will exceed $P(\theta^*, K_{\text{opt}}(\theta^*, T_s), D_{\text{U}})$ in each round with constant probability. This implies that with probability $1 - o_T(1/T)$, for every $s$, the control $u_t^{\text{safeU}}$ is used a constant fraction of the times before time $T_s$. Combining these two results, we have for all $s$ that with probability $1 - o_T(1/T)$, $\epsilon_s = \tilde{O}_T(1/\sqrt{|S_{T_s}|})$ and $|S_{T_s}| = \Omega_T(T_s)$. Therefore, we can conclude that with probability $1 - o_T(1/T)$, we have $\epsilon_s = \tilde{O}_T\left(\frac{1}{\sqrt{T_s}}\right)$.

**Regret Proof Changes**  Equipped with this tighter upper bound on $\epsilon_s$, we can bound $R_1$ (the regret of using controller $C_{K_{\text{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}$ rather than $C_{K_{\text{opt}}(\theta^*, T_s)}^{\theta^*}$) and $R_3$ (the regret of enforcing safety at every time step with controls $u_t^{\text{safeU}}$ and $u_t^{\text{safeL}}$) by $\tilde{O}_T(\sqrt{T})$ (Proposition 9 and 10, respectively). Because $\nu_T = T^{-1/4}$, $R_0$ is $\tilde{O}_T(\sqrt{T})$. Therefore, we can conclude as in the proof sketch of Theorem 2 that the regret of Algorithm 3 is upper-bounded by $R_0 + R_1 + R_2 + R_3 = \tilde{O}_T(\sqrt{T})$.

# 5    Discussion

In this paper, we have presented new results for the safety-constrained LQR problem. We conclude by discussing some possible extensions of our work and remaining open questions.

While our results focus on positional constraints, we also expect that similar results would hold for algorithms similar to Algorithms 2 and 3 when there are also constraints on the controls. While we leave the formal derivations of results for control constraints to future work, we provide a brief discussion of how the algorithm and proofs would change. With the addition of control constraints, the algorithms can no longer use $u_t^{\text{safeU}}$ or $u_t^{\text{safeL}}$ as these constraints may not satisfy the control constraints. To address this, we believe

that a slight modification to the way the algorithm chooses the controller $C_s^{\text{alg}}$ will allow the algorithms to satisfy both control and position constraints with high probability and achieve the same regret results as in Theorems 1 and 2. We propose choosing $C_s^{\text{alg}} = C_K^{\hat{\theta}_s}$, where $K$ is chosen such that it satisfies positional constraints and control constraints $\tilde{\Theta}_T(\epsilon_s)$ tighter than the actual constraints. As long as $\|\hat{\theta}_s - \theta^*\|_\infty \leq \tilde{O}_T(\epsilon_s)$, this will guarantee both types of constraints are satisfied. The main additional result that needs to be assumed is that choosing this $C_s^{\text{alg}}$ will not have significantly more regret than in the existing proofs. See Appendix I.1 for more discussion on the generalization of our results to the setting with both position and control constraints.

Our results also focus on one-dimensional LQR, but we expect that many of the same results will generalize to higher dimensions. In higher dimensions, a natural generalization of our constraints is to consider a compact safe region that is defined as the intersection of a finite number of half-planes. Therefore, the goal would be to choose controls such that the expected position stays within this safe region. We expect that the uncertainty bounds proven in this paper will generalize naturally to higher dimensions, as our bounds are based on results in Abbasi-Yadkori and Szepesvári [2011] that hold for higher dimensions. Therefore, we expect that the result of Theorem 2 will directly generalize to higher dimensions by replacing the controller $C_s^{\text{alg}}$ with $C_K^{\hat{\theta}_s}$ where $K$ is chosen as the optimal control for constraints that are $\tilde{\Theta}_T(\epsilon_s)$ tighter than the true constraints. Whether Theorem 1 generalize to higher dimensions is an open question we leave for future work, though in Appendix I.2, we discuss stylized settings in which we expect that the $\tilde{O}_T(\sqrt{T})$ regret bounds from Theorem 1 will generalize to higher dimensions.

We also note that our algorithms require knowledge of $T$ in advance, as the value of $T$ determines the length of time spent in the warm-up exploration period. We expect that similar results will hold when $T$ is not known in advance, however this would require periods of exponentially growing length that alternate exploration versus exploitation (similar to as done in, e.g. Li et al. [2021]). Because this greatly increases the complexity of the algorithm and analysis, we state and prove our results for $T$ known in advance. Finally, while we study general baselines that are more powerful than just safe linear controllers, the question of whether we can achieve $\tilde{O}_T(\sqrt{T})$ (or even $\tilde{O}_T(T^{2/3})$) regret on top of the cost of the best possible among *all* safe controllers is still open.

# Acknowledgements

# References

Yasin Abbasi-Yadkori and Csaba Szepesvári. Regret bounds for the adaptive control of linear quadratic systems. In *Proceedings of the 24th Annual Conference on Learning Theory*, pages 1–26. JMLR Workshop and Conference Proceedings, 2011.

Marc Abeille and Alessandro Lazaric. Thompson sampling for linear-quadratic control problems. In *Artificial intelligence and statistics*, pages 1246–1254. PMLR, 2017.

Archith Athrey, Othmane Mazhar, Meichen Guo, Bart De Schutter, and Shengling Shi. Regret analysis of learning-based linear quadratic gaussian control with additive exploration. In *2024 European Control Conference (ECC)*, pages 1795–1801. IEEE, 2024.

Alberto Bemporad and Manfred Morari. Robust model predictive control: A survey. In *Robustness in identification and control*, pages 207–226. Springer, 2007.

Alberto Bemporad, Manfred Morari, Vivek Dua, and Efstratios N Pistikopoulos. The explicit linear quadratic regulator for constrained systems. *Automatica*, 38(1):3–20, 2002.

Richard Cheng, Gábor Orosz, Richard M Murray, and Joel W Burdick. End-to-end safe reinforcement learning through barrier functions for safety-critical continuous control tasks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 3387–3395, 2019.

Jae Weon Choi and Young Bong Seo. Lqr design with eigenstructure assignment capability [and application to aircraft flight control]. *IEEE Transactions on Aerospace and Electronic Systems*, 35(2):700–708, 1999.

Alon Cohen, Tomer Koren, and Yishay Mansour. Learning linear-quadratic regulators efficiently with only sqrtt regret. pages 1300–1309, 2019.

Richard Combes. An extension of mcdiarmid's inequality. *arXiv preprint arXiv:1511.05240*, 2015.

Sarah Dean, Horia Mania, Nikolai Matni, Benjamin Recht, and Stephen Tu. Regret bounds for robust adaptive control of the linear quadratic regulator. *Advances in Neural Information Processing Systems*, 31, 2018.

Sarah Dean, Stephen Tu, Nikolai Matni, and Benjamin Recht. Safely learning to control the constrained linear quadratic regulator. In *2019 American Control Conference (ACC)*, pages 5582–5588. IEEE, 2019.

Mohamad Kazem Shirani Faradonbeh, Ambuj Tewari, and George Michailidis. Finite time analysis of optimal adaptive policies for linear-quadratic systems. *arXiv preprint arXiv:1711.07230*, 2017.

Mohamad Kazem Shirani Faradonbeh, Ambuj Tewari, and George Michailidis. Input perturbations for adaptive regulation and learning. *arXiv preprint arXiv:1811.04258*, 2018a.

Mohamad Kazem Shirani Faradonbeh, Ambuj Tewari, and George Michailidis. On optimality of adaptive linear-quadratic regulators. *arXiv preprint arXiv:1806.10749*, 2018b.

Charles Fefferman, Bernat Guillén Pegueroles, Clarence W Rowley, and Melanie Weber. Optimal control with learning on the fly: a toy problem. *Revista matemática iberoamericana*, 38(1):175–187, 2021.

Jaime F Fisac, Anayo K Akametalu, Melanie N Zeilinger, Shahab Kaynama, Jeremy Gillula, and Claire J Tomlin. A general safety framework for learning-based control in uncertain robotic systems. *IEEE Transactions on Automatic Control*, 64(7):2737–2752, 2018.

Nathan Fulton and André Platzer. Safe reinforcement learning via formal methods: Toward safe control through proof and learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.

Milan Ganai, Zheng Gong, Chenning Yu, Sylvia Herbert, and Sicun Gao. Iterative reachability estimation for safe reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 2024.

Kunal Garg, Songyuan Zhang, Oswin So, Charles Dawson, and Chuchu Fan. Learning safe control for multi-robot systems: Methods, verification, and open challenges. *Annual Reviews in Control*, 57:100948, 2024.

Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, and Alois Knoll. A review of safe reinforcement learning: Methods, theory and applications. *arXiv preprint arXiv:2205.10330*, 2022.

Mohammad Khosravi and Roy S Smith. Nonlinear system identification with prior knowledge on the region of attraction. *IEEE Control Systems Letters*, 5(3):1091–1096, 2020.

B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. *IEEE Transactions on Intelligent Transportation Systems*, 23(6):4909–4926, 2021.

Johannes Köhler, Elisa Andina, Raffaele Soloperto, Matthias A Müller, and Frank Allgöwer. Linear robust adaptive model predictive control: Computational complexity and conservatism. In *2019 IEEE 58th Conference on Decision and Control (CDC)*, pages 1383–1388. IEEE, 2019.

Torsten Koller, Felix Berkenkamp, Matteo Turchetta, and Andreas Krause. Learning-based model predictive control for safe exploration. In *2018 IEEE conference on decision and control (CDC)*, pages 6059–6066. IEEE, 2018.

Bruce Lee, Anders Rantzer, and Nikolai Matni. Nonasymptotic regret analysis of adaptive linear quadratic control with model misspecification. In *6th Annual Learning for Dynamics & Control Conference*, pages 980–992. PMLR, 2024.

Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep visuomotor policies. *The Journal of Machine Learning Research*, 17(1):1334–1373, 2016.

Yingying Li, Subhro Das, Jeff Shamma, and Na Li. Safe adaptive learning-based control for constrained linear quadratic regulators with regret guarantees. *arXiv preprint arXiv:2111.00411*, 2021.

Yingying Li, Tianpeng Zhang, Subhro Das, Jeff Shamma, and Na Li. Non-asymptotic system identification for linear systems with nonlinear policies. *arXiv preprint arXiv:2306.10369*, 2023.

Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.

Matthias Lorenzen, Mark Cannon, and Frank Allgöwer. Robust mpc with recursive model update. *Automatica*, 103:461–471, 2019.

Xiaonan Lu, Mark Cannon, and Denis Koksal-Rivet. Robust adaptive model predictive control: Performance and parameter estimation. *International Journal of Robust and Nonlinear Control*, 31(18):8703–8724, 2021.

Horia Mania, Stephen Tu, and Benjamin Recht. Certainty equivalence is efficient for linear quadratic control. *Advances in Neural Information Processing Systems*, 32, 2019.

Horia Mania, Michael I Jordan, and Benjamin Recht. Active learning for nonlinear system identification with guarantees. *arXiv preprint arXiv:2006.10277*, 2020.

Zahra Marvi and Bahare Kiumarsi. Safe reinforcement learning: A control barrier function optimization approach. *International Journal of Robust and Nonlinear Control*, 31(6): 1923–1940, 2021.

Ali Mesbah. Stochastic model predictive control: An overview and perspectives for future research. *IEEE Control Systems Magazine*, 36(6):30–44, 2016.

Teodor Mihai Moldovan and Pieter Abbeel. Safe exploration in markov decision processes. *arXiv preprint arXiv:1205.4810*, 2012.

Deepan Muthirayan, Jianjun Yuan, Dileep Kalathil, and Pramod P Khargonekar. Online learning for predictive control with provable regret guarantees. In *2022 IEEE 61st Conference on Decision and Control (CDC)*, pages 6666–6671. IEEE, 2022.

Frauke Oldewurtel, Colin N Jones, and Manfred Morari. A tractable approximation of chance constrained stochastic mpc based on affine disturbance feedback. In *2008 47th IEEE conference on decision and control*, pages 4731–4736. IEEE, 2008.

Samet Oymak and Necmiye Ozay. Non-asymptotic identification of lti systems from a single trajectory. In *2019 American control conference (ACC)*, pages 5655–5661. IEEE, 2019.

M Cody Priess, Richard Conway, Jongeun Choi, John M Popovich, and Clark Radcliffe. Solutions to the inverse lqr problem with application to biological systems analysis. *IEEE Transactions on control systems technology*, 23(2):770–777, 2014.

J.B. Rawlings and D.Q. Mayne. *Model Predictive Control: Theory and Design*. Nob Hill Pub., 2009. ISBN 9780975937709.

Alicia Arce Rubio, Alexandre Seuret, Yassine Ariba, and Alessio Mannisi. Optimal control strategies for load carrying drones. *Delays and Networked Control Systems*, pages 183–197, 2016.

Yahya Sattar and Samet Oymak. Non-asymptotic and accurate learning of nonlinear dynamical systems. *The Journal of Machine Learning Research*, 23(1):6248–6296, 2022.

Benjamin Schiffer and Lucas Janson. Foundations of safe online reinforcement learning in the linear quadratic regulator: $\sqrt{T}$-regret. *arXiv preprint arXiv:2504.18657*, 2025.

Karam Shabaani and Mahdi Jalili-Kharaajoo. Application of adaptive lqr with repetitive control for ups systems. In *Proceedings of 2003 IEEE Conference on Control Applications, 2003. CCA 2003.*, volume 2, pages 1124–1129. IEEE, 2003.

David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.

Max Simchowitz and Dylan Foster. Naive exploration is optimal for online lqr. In *International Conference on Machine Learning*, pages 8937–8948. PMLR, 2020.

Max Simchowitz, Horia Mania, Stephen Tu, Michael I Jordan, and Benjamin Recht. Learning without mixing: Towards a sharp analysis of linear system identification. In *Conference On Learning Theory*, pages 439–473. PMLR, 2018.

Yue Sun, Samet Oymak, and Maryam Fazel. Finite sample system identification: Optimal rates and the role of regularization. In *Learning for dynamics and control*, pages 16–25. PMLR, 2020.

Ambuj Tewari and Susan A Murphy. From ads to interventions: Contextual bandits in mobile health. *Mobile health: sensors, analytic methods, and applications*, pages 495–517, 2017.

Akifumi Wachi, Yanan Sui, Yisong Yue, and Masahiro Ono. Safe exploration and optimization of constrained mdps using gaussian processes. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.

Akifumi Wachi, Xun Shen, and Yanan Sui. A survey of constraint formulations in safe reinforcement learning. *arXiv preprint arXiv:2402.02025*, 2024.

Feicheng Wang and Lucas Janson. Exact asymptotics for linear quadratic adaptive control. *The Journal of Machine Learning Research*, 22(1):12136–12247, 2021.

Feicheng Wang and Lucas Janson. Rate-matching the regret lower-bound in the linear quadratic regulator with unknown dynamics. *arXiv preprint arXiv:2202.05799*, 2022.

Yihang Yao, Zuxin Liu, Zhepeng Cen, Jiacheng Zhu, Wenhao Yu, Tingnan Zhang, and Ding Zhao. Constraint-conditioned policy optimization for versatile safe reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 2024.

Lintao Ye, Ming Chi, Zhi-Wei Liu, and Vijay Gupta. Online actuator selection and controller design for linear quadratic regulation with unknown system model. *IEEE Transactions on Automatic Control*, 2024.

Zichen Zhao and Qianxiao Li. Adaptive sampling methods for learning dynamical systems. In *Mathematical and Scientific Machine Learning*, pages 335–350. PMLR, 2022.

Yang Zheng and Na Li. Non-asymptotic identification of linear dynamical systems using multiple trajectories. *IEEE Control Systems Letters*, 5(5):1693–1698, 2020.

Ingvar Ziemann and Henrik Sandberg. Regret lower bounds for learning linear quadratic gaussian systems. *IEEE Transactions on Automatic Control*, 2024.

# A  Notation

## A.1  Big O Notation

Throughout this paper, we use notation such as $o_T(\cdot)$, $O_T(\cdot)$, $\omega_T(\cdot)$, $\Omega_T(\cdot)$.

- $f(T) = O_T(g(T))$ if there exists $T_0$ and $M \in \mathbb{R}$ such that for $T \geq T_0$, $f(T) \leq M \cdot g(T)$.

- $f(T) = o_T(g(T))$ if for every constant $\epsilon > 0$ there exists $T_0$ such that for all $T \geq T_0$, $f(T) \leq \epsilon \cdot g(T)$.

- $f(T) = \tilde{O}_T(g(T))$ if there exists $T_0$ and $k, M \in \mathbb{R}$ such that for $T \geq T_0$, $f(T) \leq M \cdot g(T) \cdot \log^k(T)$.

Note that $\Omega_T, \omega_t$, and $\tilde{\Omega}_T$ are defined in the same way but with the inequality reversed. While this is standard notation, we want to highlight exactly how we are using this notation in our proofs. First, we note that the subscript $T$ is included to indicate that we will always be using this notation with respect to the variable $T$. Furthermore, we note that the constant $M$ that is "hidden" by the big-O notation will always be a function of known problem specification parameters, such as $q, r, \Theta, \mathcal{D}, D$. Therefore, if an expression includes an $O_T(1)$ term, this constant does not depend on any other variables in the expression. For example, suppose we state that for all $K$, $f(K) \leq O_T(\sqrt{T})$. Then this means that there exists $T_0$ and $M$ (where $M$ is a function of known problem specification parameters) such that for all $K$ and $T \geq T_0$, $f(K) \leq M \cdot \sqrt{T}$. Furthermore, we will use notation such as $f(T) = O_T(\epsilon)$ to mean that there exists $T_0$ and $M$ such that $f(T) \leq M \cdot \epsilon$ for $T \geq T_0$, where $M$ does not depend on $\epsilon$ and only depends on the problem specification parameters $\{q, r, \Theta, \mathcal{D}, D\}$. Finally, note that we will use the computer science notation of $O_T()$, in that the functions $f(T)$ and $g(t)$ will always be non-negative.

## A.2  Miscellaneous Notation

Throughout the proofs, any inequalities or equations involving random variables will represent inequality or equality almost surely unless otherwise stated. Throughout the paper, we will use the notation $\{x_i\}_{i=1}^n$ to represent the unordered but indexed set of $x_1, x_2, ..., x_n$.

## A.3 Problem Specifications

The notation below will be used throughout the appendix, however the variables may depend on the algorithm being studied within a section. For example, the event $E$ is defined slightly differently for each of the two algorithms, and therefore the reader should note which algorithm each section addresses. The notation never changes within a single section.

- $q, r$ : coefficients for the cost at time $t$ of $qx_t^2 + ru_t^2$.

- $W = \{w_t\}_{t=0}^{T-1}$ : The noise random variables for the $T$-length trajectory.

- $\mathcal{D}$ : Distribution of $w_t$ with CDF $F_{\mathcal{D}}$ and pdf upper bound $B_P$

- $\Theta = [\underline{a}, \bar{a}] \times [\underline{b}, \bar{b}]$ : Given set of dynamics s.t. $\theta^* \in \Theta$ (size$(\Theta) = \min(\bar{a} - \underline{a}, \bar{b} - \underline{b})$)

- $\theta^* = (a^*, b^*)$ : The true (unknown) dynamics.

- $C^{\text{init}}$ : The initial safe controller satisfying Assumption 1.

- $D = (D_{\text{L}}, D_{\text{U}})$ : the expected-position boundary for the safety constraint.

- A set of controls $\{u_t\}$ are safe for dynamics $\{\theta_t\}$ if for all $t$, $D_{\text{L}} \le a_t x_t + b_t u_t \le D_{\text{U}}$.

- $H_t = (x_0, u_0, x_1, u_1, ..., u_{t-1}, x_t)$ and $\mathcal{F}_t = \sigma(H_t)$.

- $J(\theta, C, T, x, W)$ : The random variable cost of using controller $C$ starting at position $x_0 = x$ for $T$ time steps under dynamics $\theta$ with noise random variables $W$.

- $J^*(\theta, C, T) = J^*(\theta, C, T, 0) = \mathbb{E}[J(\theta, C, T, x, W) \mid \theta, C, T, x]$ and $J^*(\theta, C, T) = J^*(\theta, C, T, 0)$.

- $J^*(\theta, C) = J^*(\theta, C, 0) = \lim_{T \to \infty} J^*(\theta, C, T, 0)$.

- $\mathcal{C}^\theta = \{C_K^\theta\}_{K \in [K_{\text{L}}^\theta, K_{\text{U}}^\theta]}$ : a class of controllers that are safe for dynamics $\theta$

- $K_{\text{opt}}(\theta, T)$ : The $K$ that maximizes $J^*(\theta, C_K^\theta, T, 0)$ for $K \in [K_{\text{L}}^\theta, K_{\text{U}}^\theta]$.

- $K_{\text{opt}}(\theta)$ : The $K$ that maximizes $J^*(\theta, C_K^\theta)$ for $K \in [K_{\text{L}}^\theta, K_{\text{U}}^\theta]$.

- $C_K^{\text{unc}}$ : The unconstrained linear controller with parameter $K$, i.e. $C_K^{\text{unc}}(x) = -Kx$.

- $F_{\text{opt}}(\theta)$ : The $K$ that maximizes $J^*(\theta, C_K^{\text{unc}})$.

## A.4 Algorithm Notation

- $\nu_T$ : Algorithm specific parameter that is either $T^{-1/4}$ or $T^{-1/3}$.

- $s_e$ : The number of the last round of the safe exploitation phase.

- $T_s = \frac{2^s}{\nu_T^2}$ : The length **and** starting time of round $s$ of the safe exploitation phase.

- $\epsilon_s$ : Uncertainty bound for $\theta^*$ used throughout the algorithm.

- $\hat{\theta}_s$ : An estimate of $\theta^*$ that is with high probability within $\epsilon_s$ distance of $\theta^*$

- $u_t^{\text{safeU}}$ : Largest $u$ such that $\max\limits_{\|\theta-\hat{\theta}_s\|_\infty \leq \epsilon_s} ax_t + bu \leq D_{\text{U}}$

- $u_t^{\text{safeL}}$ : Smallest $u$ such that $\max\limits_{\|\theta-\hat{\theta}_s\|_\infty \leq \epsilon_s} ax_t + bu \geq D_{\text{L}}$.

- $C_s^{\text{alg}}(x_t)$ : the controller that the algorithm uses in round $s$ of the safe exploitation phase with additional safety modifications, i.e. the algorithm in round $s$ of the safe exploitation phase uses control $u_t = \max\left(\min\left(C_s^{\text{alg}}(x_t), u_t^{\text{safeU}}\right), u_t^{\text{safeL}}\right)$.

- $C^{\text{alg}}$ : Controller of the corresponding algorithm as described in the previous point.

- $P(\theta, K, z)$ : See Assumption 9.

## A.5 Proof Notation

- $W_s = \{w_i\}_{i=T_s}^{T_{s+1}-1}$ : Noise random variables in the round $s$ of safe exploitation phase.

- $\left(C_{K^*}^{\theta^*}, \{C_{K_s^*}^{\theta^*}\}_{s=0}^{s_e}\right)$ : The expected cost minimizing set of controllers to use if the controller $C_{K^*}^{\theta^*}$ is used for the first $T_0$ steps and for time $t \geq T_0$, the controller used is $C_{K_s}^{\theta^*}$, where $s = \lfloor \log_2(t\nu_T^2) \rfloor$. The sequence $(x_0^*, x_1^*, ...)$ are the corresponding positions of using these controllers.

- $(x_0', x_1', ...)$ and $(u_0', u_1', ...)$: Unless otherwise specified, these are the positions and controls of the algorithm being discussed in the current proof.

- $(\hat{x}_{T_0}, \hat{x}_{T_0+1}, ...)$ : Unless otherwise defined in the theorem/lemma statement, $\hat{x}_{T_0}, \hat{x}_{T_0+1}, ...$ is the sequence of positions if the control at each time $t \geq T_0$ is $C_{K_{\text{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}(x_t)$ for $s = \lfloor \log_2(t\nu_T^2) \rfloor$ and starting at $\hat{x}_{T_0} = x_{T_0}'$.

- $E_{\text{safe}} = \{\forall t < T : D_{\text{L}} \leq a^*x_t' + b^*u_t' \leq D_{\text{U}}\}$ : Event that all controls are safe

- $E_1 = \{\forall t < T : |w_t| \leq \log^2(T)\}$ : Event that all noise has magnitude less than $\log^2(T)$

- $E_0 = \{\forall s \leq s_e : \|\theta^* - \hat{\theta}_s\|_\infty \leq \epsilon_s\}$ : Event that all estimates of $\theta^*$ are within $\epsilon_s$ of $\theta^*$.

- $E_2 = E_0 \bigcap \{\max_{s\in[0:s_e]} \epsilon_s \leq \tilde{O}_T(\nu_T)\}$.

- $E_2^s = \{\|\hat{\theta}_s - \theta^*\|_\infty \leq \epsilon_s \leq c_T \cdot \nu_T\}$, where $c_T$ is the coefficient in the $\tilde{O}_T(\nu_T)$ of the definition of event $E_2$.

- $E = E_{\text{safe}} \cap E_1 \cap E_2$

- $B_x = \log^3(T)$ : Used throughout the appendix to simplify notation.

# B    Additional Related Work

The constrained LQR problem is closely related to the problem of model predict control (MPC) with constraints. For example, there is a large body of work on robust model predictive control with known dynamics [Bemporad and Morari, 2007]. This is further extended to MPC with model uncertainties in robust adaptive MPC (RAMPC) in works such as Köhler et al. [2019], Lu et al. [2021]. There have also been significant work on stochastic MPC with soft constraints, for example Mesbah [2016], Oldewurtel et al. [2008], which are closely related to the expected position constraints we use in this paper. In the context of constrained LQR with no noise, Bemporad et al. [2002] derive the optimal controller as a piece-wise affine function. In a different MPC setting with deterministic dynamics and noisy observations, Muthirayan et al. [2022] provide an algorithm that also achieves $O(T^{2/3})$ regret. Learning based MPC using an initial safe controller was also studied in Koller et al. [2018]. MPC results on learning constraints include e.g. Lorenzen et al. [2019], Köhler et al. [2019]. While these works provide algorithms to solve constrained optimization problems such as LQR, these works do not compare the asymptotic performance of their results to the optimal algorithm. In contrast, our work studies a similar problem but focuses on algorithmic regret analysis from an RL perspective, comparing our algorithm to some baseline representation of the "best" algorithm.

The results in this paper are also closely related to general system identification, the idea of being able to (in any way) asymptotically estimate the unknown dynamics. There have been multiple works in this area including Simchowitz et al. [2018], Zhao and Li [2022], Mania et al. [2020]. A recent work closely related to the results of this paper is Li et al. [2023], which describes learning rates for non-linear controllers in a similar setting. The results in Li et al. [2023], however, require i.i.d. noise excitation in every step, while our uncertainty bounds after the warm-up phase actually require no such excitation. These works are most similar to our work in that our results rely on identifying the system dynamics to a high accuracy. However our focus is not simply on learning the system, but also on achieving provably low regret results. The new uncertainty bounds we use to achieve our results also apply to nonlinear controllers as in Li et al. [2023], but our uncertainty bounds apply specifically to the setting with safety constraints.

# C    Proof of Theorem 2

Before proving Theorem 2, we extend Definition 1 to account for time-dependent dynamics.

**Definition 2.** *A control $u_t$ and position $x_t$ are safe for dynamics $\theta_t$ if*

$$D_{\mathrm{L}} \leq a_t x_t + b_t u_t \leq D_{\mathrm{U}}.$$

*Similarly, a (possibly time-dependent) controller $C_t$ is safe for $T$ steps for dynamics $\{\theta_t\}$ if when the dynamics at time $t$ is $\theta_t$, the sequence of controls $C_0(H_0), C_1(H_1), ..., C_{T-1}(H_{T-1})$ and the resulting positions $x_0, ..., x_{T-1}$ are safe for dynamics $\theta_t$ at all times $t$.*

Note that in general, a controller being safe is a random event.

Theorem 2 makes two claims: the first is that Algorithm 2 is safe for dynamics $\theta^*$ for all $T$ steps with high probability and the second bounds with high probability the regret of Algorithm 2. In Appendix C.1 we will prove the result about the safety of Algorithm 2 and in Appendix C.2 we will prove the result about the regret of Algorithm 2.

## C.1   Proof of Safety of Algorithm 2

**Lemma 1.** *Under Assumptions 1–8 , Algorithm 2 is safe for $T$ steps for dynamics $\theta^*$ with probability $1 - o_T(1/T^2)$.*

*proof.* We will first analyze the warm-up exploration phase (the first loop in Algorithm 2 in Lines 2–4). If the control at time $t-1$ was safe for dynamics $\theta^*$ as in Definition 2, then with probability at least $1 - O_T(\frac{1}{T^4})$, the next position satisfies

$$x_t \in \left[ D_{\mathrm{L}} - F_{\mathcal{D}}^{-1}(1 - \frac{1}{T^4}), D_{\mathrm{U}} + F_{\mathcal{D}}^{-1}(1 - \frac{1}{T^4}) \right].$$

By Assumption 2 on the controller $C^{\mathrm{init}}$, $D_{\mathrm{L}} + \frac{b^*}{\log(T)} \leq a^*x + b^*C^{\mathrm{init}}(x) \leq D_{\mathrm{U}} - \frac{b^*}{\log(T)}$ for all $x \in \left[ D_{\mathrm{L}} - F_{\mathcal{D}}^{-1}(1 - \frac{1}{T^4}), D_{\mathrm{U}} + F_{\mathcal{D}}^{-1}(1 - \frac{1}{T^4}) \right]$. In Lines 2–4 of Algorithm 2 the control is $C^{\mathrm{init}}(x_t) + \frac{\phi_t}{\log(T)}$ and $|\phi_t| = 1$. Therefore, if at time $t-1$ the algorithm's control was safe, then with probability $1 - O_T\left(\frac{1}{T^4}\right)$ the control at time $t$ will satisfy $D_{\mathrm{L}} \leq a^*x_t + b^*u_t \leq D_{\mathrm{U}}$ and be safe. Furthermore, at time 0, the position is $x_0 = 0$, therefore the first control is safe. Using this as a base case in a proof by induction with a union bound over all $1/\nu_T^2$ time steps $t$ in this loop, with probability $1 - O_T(1/T^3)$, the first $1/\nu_T^2$ steps will be safe for dynamics $\theta^*$.

Now we will analyze the second loop in Algorithm 2 (Lines 5–13). Define $s_e = \log_2(T\nu_T^2) - 1$. Define the event $E_0$ as

$$E_0 = \left\{ \forall s \leq s_e : \|\theta^* - \hat{\theta}_s\|_\infty \leq \epsilon_s \right\}. \tag{15}$$

These $\epsilon_s$ are less than the right hand side of the equation in Lemma 23, and therefore by Lemma 23, under Assumptions 3 and 1,

$$\mathbb{P}(E_0) \geq 1 - o_T(1/T^2). \tag{16}$$

Informally, the next event we define is the combination of event $E_0$ and the event that the $\epsilon_s$ (defined in Line 9 of Algorithm 2) are decreasing at a sufficiently fast rate, which we will prove in Lemma 2. Define

$$E_2 = E_0 \bigcap \left\{ \max_{s \in [0:s_e]} \epsilon_s \leq \tilde{O}_T(\nu_T) \right\}. \tag{17}$$

**Lemma 2.** *Under Assumptions 1–8, with probability $1 - o_T(1/T^2)$*

$$\max_{s \in [0:s_e]} \epsilon_s \leq \tilde{O}_T(\nu_T).$$

The proof of Lemma 2 can be found in Appendix G.2. Combining Lemma 23 and Lemma 2 with a union bound gives that

$$\mathbb{P}(E_2) \geq 1 - o_T(1/T^2). \tag{18}$$

Define the event $E_1$ as

$$E_1 = \left\{ \forall t < T : |w_t| \leq \log^2(T) \right\}. \tag{19}$$

By Assumption 3, the noise is sub-Gaussian, and therefore there exists a constant $\alpha$ such that for any $t$ and $x$, $\mathbb{P}(w_t \geq x) \leq 2\exp(-x^2/\alpha)$. Taking $x = \log^2(T)$ and a union bound over all $w_t$, we have that

$$\mathbb{P}(E_1) \geq 1 - \sum_{t=0}^{T-1} 2\exp\left(-\log^4(T)/\alpha\right) = 1 - o_T\left(\frac{1}{T^{\log(T)}}\right). \tag{20}$$

We need one last lemma before concluding the proof.

**Lemma 3.** *Under Assumptions 1–8, conditional on $E_1 \cap E_2$ and for sufficiently large $T$, if $u_{T_0-1}$ is safe for dynamics $\theta^*$, then for all $t \in [T_0, T]$,*

$$u_t^{\text{safeL}} \leq u_t^{\text{safeU}}.$$

The proof of Lemma 3 can be found in Appendix E.1.

Under event $E_0$, $\hat{\theta}_s$ satisfies $\|\theta^* - \hat{\theta}_s\|_\infty \leq \epsilon_s$ for all $s \in [0 : s_e]$ (which recall are the $s$ in the second for loop of Algorithm 2). Therefore, by the choice of $u_t^{\text{safeU}}$ and $u_t^{\text{safeL}}$ in Lines 11 and 12, it must be the case that $a^* x_t + b^* u_t^{\text{safeU}} \leq D_U$ and $a^* x_t + b^* u_t^{\text{safeL}} \geq D_L$. By the choice of $u_t$ in Line 13 of Algorithm 2, if $u_t^{\text{safeL}} \leq u_t^{\text{safeU}}$ then $u_t^{\text{safeL}} \leq u_t \leq u_t^{\text{safeU}}$. This implies that

$$D_L \leq a^* x_t + b^* u_t \leq D_U. \tag{21}$$

Therefore, by Lemma 3, under $E_1 \cap E_2 \cap \{u_{T_0-1}$ is safe for dynamics $\theta^*\}$, all controls used in the second for loop (Lines 5–13) in Algorithm 2 are safe for dynamics $\theta^*$. By a union bound combining Equations (18) and (20) and the first paragraph of this proof, we have that

$$\mathbb{P}(E_1 \cap E_2 \cap \{u_{T_0-1} \text{ is safe for dynamics } \theta^*\}) = 1 - o_T(1/T^2).$$

Because all of the steps in Algorithm 2 are part of either the first or second loop, and the first loop steps are safe for dynamics $\theta^*$ with probability $1 - o_T(1/T^2)$ and the second loop steps are safe for dynamics $\theta^*$ with probability $1 - o_T(1/T^2)$, a union bound gives that the overall algorithm is safe for dynamics $\theta^*$ with probability $1 - o_T(1/T^2)$. $\qquad\square$

## C.2 Proof of Regret Bound of Algorithm 2

*proof.* Define the event $E_{\text{safe}}$ as the event that the controls used by the algorithm are safe at all times. If $x'_0, x'_1, \ldots$ and $u'_0, u'_1, \ldots$ are respectively the positions and controls of the algorithm, we have that

$$E_{\text{safe}} = \left\{ \forall t < T : D_L \leq a^* x'_t + b^* u'_t \leq D_U \right\}, \tag{22}$$

and by Lemma 1 we have that $\mathbb{P}(E_{\text{safe}}) = 1 - o_T(1/T^2)$. Now, define the event $E$ as

$$E = E_{\text{safe}} \cap E_1 \cap E_2. \tag{23}$$

A union bound combining Equations (20) and (18) gives that

$$\mathbb{P}(E) = \mathbb{P}(E_{\text{safe}} \cap E_1 \cap E_2) \geq 1 - o_T(1/T^2). \tag{24}$$

The rest of the proof of Theorem 2 will focus on proving that the regret of Algorithm 2 is $\tilde{O}_T(T^{2/3})$ with conditional probability at least $1 - o_T(1/T)$ given $E$. Let $C^{\text{alg}}$ be the (time-dependent) controller of Algorithm 2. Then the total cost of using Algorithm 2 is $T \cdot J(\theta^*, C^{\text{alg}}, T, 0, W)$, and the regret we are trying to bound is (as in Equation (7) using the notation $K_{\text{opt}}$ from Assumption 6),

$$T \cdot J(\theta^*, C^{\text{alg}}, T, 0, W) - T \cdot \bar{J}(\theta^*, C^{\theta^*}_{K_{\text{opt}}(\theta^*, T)}, T). \tag{25}$$

Define $W_s$ as the noise random variables from time $T_s$ to $T_{s+1} - 1$, so

$$W_s = \{w_i\}_{i=T_s}^{T_{s+1}-1}. \tag{26}$$

For any tuple $(K, \{K_s\}_{0 \leq s \leq s_e})$ where $K, K_s \in (K_L^{\theta^*}, K_U^{\theta^*})$, define $x_0^{(K, \{K_s\}_{0 \leq s \leq s_e})}, x_1^{(K, \{K_s\}_{0 \leq s \leq s_e})}, \dots$ as the random variable sequence of positions that result from starting at $x_0 = 0$ and using the controller that at each time $t < T_0$ uses controller $C_K^{\theta^*}$ and at each time $t \geq T_0$ uses the controller $C_{K_s}^{\theta^*}$, where $s = \lfloor \log_2 (t\nu_T^2) \rfloor$. Define $(K^*, \{K_s^*\}_{0 \leq s \leq s_e})$ as follows:

$$(K^*, \{K_s^*\}_{0 \leq s \leq s_e})$$
$$= \underset{(K, \{K_s\}_{0 \leq s \leq s_e})}{\arg\min} \mathbb{E}\left[ \frac{1}{\nu_T^2} J\left( \theta^*, C_K^{\theta^*}, \frac{1}{\nu_T^2}, 0, \{w_t\}_{t=0}^{T_0-1} \right) + \sum_{s=0}^{s_e} T_s J(\theta^*, C_{K_s}^{\theta^*}, T_s, x_{T_s}^{(K, \{K_s\}_{0 \leq s \leq s_e})}, W_s) \right].$$

Here the expectation is taken over both $w_t$ and $W_s$ (and recall that $x_{T_s}$ is a deterministic function of the $w_t$ and $W_s$ because $C_K^\theta$ is non-random for all $K, \theta$). We then define $x_0^*, x_1^*, \dots$ as the random variable sequence of positions such that $x_t^* = x_t^{(K^*, \{K_s^*\}_{0 \leq s \leq s_e})}$. By construction, we could choose $K, K_s = K_{\text{opt}}(\theta^*, T)$ for every $s$, and therefore it must be the case that

$$\mathbb{E}\left[ \frac{1}{\nu_T^2} J\left( \theta^*, C_{K^*}^{\theta^*}, \frac{1}{\nu_T^2}, 0, \{w_t\}_{t=0}^{T_0-1} \right) + \sum_{s=0}^{s_e} T_s J(\theta^*, C_{K_s^*}^{\theta^*}, T_s, x_{T_s}^*, W_s) \right] \leq T \cdot \bar{J}\left( \theta^*, C^{\theta^*}_{K_{\text{opt}}(\theta^*, T)}, T \right).$$

Therefore, upper bounding the cost of Algorithm 2 minus the cost of using $K^*$ for $T_0$ steps and then using the sequence of controllers $\{C_{K_s^*}^{\theta^*}\}$ each for $T_s$ steps is sufficient for upper bounding the regret in Equation (25). Now we will bound

$$T \cdot J(\theta^*, C^{\text{alg}}, T, 0, W) - \mathbb{E}\left[ \frac{1}{\nu_T^2} J\left( \theta^*, C_{K^*}^{\theta^*}, \frac{1}{\nu_T^2}, 0, \{w_t\}_{t=0}^{T_0-1} \right) + \sum_{s=0}^{s_e} T_s J(\theta^*, C_{K_s^*}^{\theta^*}, T_s, x_{T_s}^*, W_s) \right].$$
$$\tag{27}$$

Note that we will upper bound the cost in terms of the parameter $\nu_T = T^{-1/3}$ in Line 1. In order to bound the quantity in Equation (27), we will break this component of regret

into four sources: the regret from the warm-up period (Lines 2–4), the regret from using the estimates $\hat{\theta}_s$ instead of using $\theta^*$, the regret induced by the randomness of the trajectory, and the regret from enforcing safety.

The first source of regret is the regret incurred in the warm-up period of Algorithm 2 (Lines 2–4). Recall that $C_s^{\text{alg}}$ is the controller used in Algorithm 2 in the $s$ iteration of the second for loop. We will use Proposition 3 to bound the cost incurred during the warm-up period.

**Proposition 3** (Regret from Warm-up Period). *Define $x_0', x_1', \ldots$ as the sequence of random variables that are the positions of the controller $C^{\text{alg}}$ defined in Algorithm 2. Define $R_0$ as the cost of the first $1/\nu_T^2$ steps, i.e.*

$$R_0 = T \cdot J(\theta^*, C^{\text{alg}}, T, 0, W) - \sum_{s=0}^{s_e} T_s \cdot J(\theta^*, C_s^{\text{alg}}, T_s, x_{T_s}', W_s). \tag{28}$$

*Then under Assumptions 1–8 and conditional on event $E$,*

$$R_0 \overset{a.s.}{\leq} \tilde{O}_T\left(\frac{1}{\nu_T^2}\right).$$

The proof of Propposition 3 can be found in Appendix D.1. The second source of regret in Equation (27) is that Algorithm 2 uses a controller $C_{K_{\text{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}$ instead of the controller $C_{K_s^*}^{\theta^*}$. This source of regret (denoted $R_1$) can be interpreted as the "estimation cost" of using the estimated controller instead of the optimal controller, but without enforcing safety. We will use Proposition 4 to bound this source of regret.

**Proposition 4** (Regret from Non-optimal Controller). *Define $R_1$ as*

$$R_1 := \sum_{s=0}^{s_e} \mathbb{E}\left[T_s J(\theta^*, C_{K_{\text{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}, T_s, 0, W_s) \,\Big|\, \hat{\theta}_s\right] - \mathbb{E}\left[\sum_{s=0}^{s_e} T_s J(\theta^*, C_{K_s^*}^{\theta^*}, T_s, x_{T_s}^*, W_s)\right].$$

*Note that $W_s$ is independent of $\hat{\theta}_s$ by construction. Then under Assumptions 1–8 and conditional on event $E_2$,*

$$R_1 \overset{a.s.}{\leq} \tilde{O}_T(T\nu_T). \tag{29}$$

The proof of Proposition 4 can be found in Appendix D.2. It may appear odd that the starting positions of the two terms do not match in the definition of $R_1$ (or in the definition of $R_2$ below), but we do account for this difference in the proofs of Propositions 4 and 5. The third source of regret (which we will denote $R_2$) comes from the fact that in Equation (27) we are comparing the random variable $T \cdot J(\theta^*, C^{\text{alg}}, T, 0, W)$ to an expectation. In order to show that this source of regret is small, we need to show a concentration inequality for the cost of repeatedly using controllers of the form $C_{K_{\text{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}$, which we do in Proposition 5.

**Proposition 5** (Regret from Randomness). *Define $\hat{x}_{T_0}, \hat{x}_{T_0+1}, \ldots$ as the sequence of random variables representing the sequence of positions if the control at each time $t \geq T_0$ is $C_{K_{\text{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}(x_t)$ for $s = \lfloor \log_2(t\nu_T^2)\rfloor$ and starting at $\hat{x}_{T_0} = x_{T_0}'$. Define $R_2$ as*

$$R_2 := \sum_{s=0}^{s_e} T_s J(\theta^*, C_{K_{\text{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}, T_s, \hat{x}_{T_s}, W_s) - \sum_{s=0}^{s_e} \mathbb{E}\left[T_s J(\theta^*, C_{K_{\text{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}, T_s, 0, W_s) \,\Big|\, \hat{\theta}_s\right].$$

*Then with conditional probability $1 - o_T(1/T)$ given event $E$,*

$$R_2 \leq \tilde{O}_T(\sqrt{T}). \tag{30}$$

The proof of Proposition 5 can be found in Appendix D.3. The final source of regret in Equation (27) is the extra cost incurred by enforcing safety in Algorithm 2 (Line 13) rather than using the control given by $C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}$. Each time we enforce safety we potentially incur an extra cost, but Proposition 6 bounds this extra cost.

**Proposition 6** (Regret from Enforcing Safety). *Define $\hat{x}_{T_0}, \hat{x}_{T_0+1}, \dots$ as the sequence of random variables representing the sequence of positions if the control at each time $t \geq T_0$ is $C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x_t)$ for $s = \lfloor \log_2(t\nu_T^2) \rfloor$ and starting at $\hat{x}_{T_0} = x'_{T_0}$. Define $R_3$ as (the random variable)*

$$R_3 := \sum_{s=0}^{s_e} T_s J(\theta^*, C_s^{\text{alg}}, T_s, x'_{T_s}, W_s) - \sum_{s=0}^{s_e} T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_s, \hat{x}_{T_s}, W_s).$$

*Then under Assumptions 1–8, with conditional probability $1 - o_T(1/T)$ given event $E$,*

$$R_3 \leq \tilde{O}_T(\nu_T T).$$

The proof of Proposition 6 can be found in Appendix D.4. Now we are ready to combine all of the sources of regret. To summarize, we have bounded and broken down the regret into

$$T \cdot J(\theta^*, C^{\text{alg}}, T, 0, W) - T \cdot \bar{J}(\theta^*, C^{\theta^*}_{K_{\text{opt}}(\theta^*, T)}, T)$$

$$\leq T \cdot J(\theta^*, C^{\text{alg}}, T, 0, W) - \mathbb{E}\left[\frac{1}{\nu_T^2} J\left(\theta^*, C^{\theta^*}_{K^*}, \frac{1}{\nu_T^2}, 0, \{w_t\}_{t=0}^{1/\nu_T^2-1}\right) + \sum_{s=0}^{s_e} T_s J(\theta^*, C^{\theta^*}_{K_s^*}, T_s, x^*_{T_s}, W_s)\right]$$

$$\leq T \cdot J(\theta^*, C^{\text{alg}}, T, 0, W) - \mathbb{E}\left[\sum_{s=0}^{s_e} T_s \bar{J}(\theta^*, C^{\theta^*}_{K_s^*}, T_s, x^*_{T_s}, W_s)\right]$$

$$= \underbrace{\sum_{s=0}^{s_e} \mathbb{E}\left[T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \,\Big|\, \hat{\theta}_s\right] - \mathbb{E}\left[\sum_{s=0}^{s_e} T_s J(\theta^*, C^{\theta^*}_{K_s^*}, T_s, x^*_{T_s}, W_s)\right]}_{R_1}$$

$$+ \underbrace{\sum_{s=0}^{s_e} T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_s, \hat{x}_{T_s}, W_s) - \sum_{s=0}^{s_e} \mathbb{E}\left[T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \,\Big|\, \hat{\theta}_s\right]}_{R_2}$$

$$+ \underbrace{\sum_{s=0}^{s_e} T_s J(\theta^*, C_s^{\text{alg}}, T_s, x'_{T_s}, W_s) - \sum_{s=0}^{s_e} T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_s, \hat{x}_{T_s}, W_s)}_{R_3}$$

$$+ \underbrace{T \cdot J(\theta^*, C^{\text{alg}}, T, 0, W) - \sum_{s=0}^{s_e} T_s J(\theta^*, C_s^{\text{alg}}, T_s, x'_{T_s}, W_s)}_{R_0}. \tag{31}$$

Now we will use Propositions 3, 4, 5, and 6 to bound the above quantity. Conditional on event $E$, Proposition 3 and Proposition 4 respectively imply that $R_0 \leq \tilde{O}_T(1/\nu_T^2)$ and $R_1 \leq \tilde{O}_T(\nu_T T)$. Proposition 5 and Proposition 6 respectively imply that conditional on event $E$ with conditional probability $1 - o_T(1/T)$, $R_2 \leq \tilde{O}_T(\sqrt{T})$ and $R_3 \leq \tilde{O}_T(\nu_T T)$. Therefore, applying a union bound gives that the bounds on $R_0, R_1, R_2, R_3$ all hold conditional on event $E$ with probability $1 - o_T(1/T)$. Putting these bounds into Equation (31), we have that conditional on event $E$ with probability $1 - o_T(1/T)$,

$$T \cdot J(\theta^*, C^{\text{alg}}, T, 0, W) - T \cdot \bar{J}(\theta^*, C^{\theta^*}_{K_{\text{opt}}(\theta^*, T)}, T) \leq R_1 + R_2 + R_3 + R_0 \leq \tilde{O}_T\left(\sqrt{T} + \frac{1}{\nu_T^2} + T\nu_T\right).$$

Choosing $\nu_T = T^{-1/3}$ (as in Algorithm 2) will minimize this regret upper bound giving a total regret upper bound of $\tilde{O}_T(T^{2/3})$. Because the probability of event $E$ is $1 - o_T(1/T)$, by a union bound the regret bound holds with unconditional probability $1 - o_T(1/T)$. $\qquad \square$

# D   Proofs of Propositions from Appendix C

## D.1   Proof of Proposition 3 (Regret of Warm-up)

*proof.* To bound the cost of the warm-up phase, we need the following lemma. Informally, Lemma 4 shows that when the noise is relatively small and the controller is "close" to being safe with respect to dynamics $\theta^*$, the position stays relatively small. Note that in this lemma we define $B_x := \log^3(T)$, which we will use throughout the proofs in the rest of the appendices.

**Lemma 4.** *Let $|x_0| \leq 4\log^2(T)$. Suppose for all $t < T$, the control used by controller $C_t$ at time $t$ is safe for fixed dynamics $\theta_t$ and for all $t \leq T$,*

$$\|\theta^* - \theta_t\|_\infty \leq \frac{1}{\log(T)}. \tag{32}$$

*Then under Assumptions 1–8, for sufficiently large $T$ and conditioned on event $E_1$, using this controller $C_t$ with dynamics $\theta^*$ for $T$ steps starting at $x_0$ will give positions $(x_0, ..., x_T)$ and controls $(u_0, ..., u_{T-1})$ satisfying the following equations.*

$$|x_t| \overset{a.s.}{\leq} 4\log^2(T) < \log^3(T) := B_x \tag{33}$$

$$|u_t| \overset{a.s.}{\leq} O_T(\log^2(T)) < \log^3(T) := B_x. \tag{34}$$

*Furthermore, if $x_0$ and the controller $C_t$ are deterministic, then the positions $(x_0, ..., x_T)$ and controls $(u_0, ..., u_{T-1})$ satisfy*

$$\mathbb{E}[|x_t|] \leq 4\log^2(T) < \log^3(T) := B_x \tag{35}$$

$$\mathbb{E}[|u_t|] \leq O_T(\log^2(T)) < \log^3(T) := B_x. \tag{36}$$

33

The proof of Lemma 4 can be found in Appendix E.2.

Now we will use this lemma to bound the total cost of the warm-up phase of the algorithm. The controller for the first $1/\nu_T^2$ steps is safe for dynamics $\theta^*$ under event $E$ as shown in Lemma 1. This means by Lemma 4, conditional on event $E$, the position and controls during this warm-up period are both bounded in magnitude by $B_x$ (defined in Lemma 4) almost surely for sufficiently large $T$. Because the cost at time $t$ is $qx_t^2 + ru_t^2$, this implies that the total cost of the first $1/\nu_T^2$ steps is upper bounded by $O_T((q+r)\frac{B_x^2}{\nu_T^2}) = \tilde{O}_T(1/\nu_T^2)$. $\qquad\square$

## D.2 Proof of Proposition 4 (Regret of Non-optimal Controller)

*proof.* First, we will use Lemma 5 to rewrite the expression in Proposition 4 in a form amenable to Assumption 7.

**Lemma 5.** *Under Assumptions 1–8 , for every $s \in [0 : s_e]$ the following hold.*

$$\left| \mathbb{E}\left[ T_s J(\theta^*, C_{K_s^*}^{\theta^*}, T_s, x_{T_s}^*, W_s) \right] - \mathbb{E}\left[ T_s J(\theta^*, C_{K_s^*}^{\theta^*}, T_s, 0, W_s) \right] \right| \leq \tilde{O}_T(1) \qquad (37)$$

The proof of Lemma 5 can be found in Appendix E.3. By Lemma 2, there exists a $c_T = \tilde{O}_T(1)$ such that under event $E_2$, $\max_s \epsilon_s \leq c_T \cdot \nu_T$. For $s \in [0 : s_e]$, define

$$E_2^s = \left\{ \|\hat{\theta}_s - \theta^*\|_\infty \leq \epsilon_s \leq c_T \cdot \nu_T \right\}. \qquad (38)$$

Informally, the event $E_2^s$ is the event that the bounds in event $E_2$ hold at time $s$. Note that because $E_2^s \subseteq E_2$, by Equation (18),

$$\mathbb{P}(E_2^s) \geq \mathbb{P}(E_2) \geq 1 - o_T(1/T^2). \qquad (39)$$

We will also use the following application of Assumption 7 that holds under event $E_2^s$. Conditional on event $E_2^s$,

$$\left| \mathbb{E}\left[ T_s J(\theta^*, C_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}, T_s, 0, W_s) - T_s J(\theta^*, C_{K_s^*}^{\theta^*}, T_s, 0, W_s) \,\middle|\, \hat{\theta}_s \right] \right|$$
$$= \left| T_s \bar{J}(\theta^*, C_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}, T_s) - T_s \bar{J}(\theta^*, C_{K_s^*}^{\theta^*}, T_s) \right|$$
$$\leq \tilde{O}_T\left( T_s \epsilon_s + \frac{T_s}{T^2} \right). \qquad \text{Assumption 7} \qquad (40)$$

We can now use the triangle inequality with Equation (37) to rewrite the left side of Equation

(29) and apply Equation (40). Formally, conditional on event $E_2$,

$$\sum_{s=0}^{s_e} \mathbb{E}\left[T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \,\Big|\, \hat{\theta}_s\right] - \mathbb{E}\left[\sum_{s=0}^{s_e} T_s J(\theta^*, C^{\theta^*}_{K^*_s}, T_s, x^*_{T_s}, W_s)\right]$$

$$= \sum_{s=0}^{s_e} \mathbb{E}\left[T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \,\Big|\, \hat{\theta}_s\right] - \sum_{s=0}^{s_e} \mathbb{E}\left[T_s J(\theta^*, C^{\theta^*}_{K^*_s}, T_s, x^*_{T_s}, W_s)\right]$$

$$\leq \tilde{O}_T(1) + \sum_{s=0}^{s_e} \mathbb{E}\left[T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \,\Big|\, \hat{\theta}_s\right] - \sum_{s=0}^{s_e} \mathbb{E}\left[T_s J(\theta^*, C^{\theta^*}_{K^*_s}, T_s, 0, W_s)\right] \quad \text{By Equation (37)}$$

$$= \tilde{O}_T(1) + \sum_{s=0}^{s_e} \mathbb{E}\left[T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) - T_s J(\theta^*, C^{\theta^*}_{K^*_s}, T_s, 0, W_s) \,\Big|\, \hat{\theta}_s\right]$$

$$\leq \tilde{O}_T(1) + \sum_{s=0}^{s_e} \left|\mathbb{E}\left[T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) - T_s J(\theta^*, C^{\theta^*}_{K^*_s}, T_s, 0, W_s) \,\Big|\, \hat{\theta}_s\right]\right|$$

$$\leq \tilde{O}_T(1) + \tilde{O}_T\left(\sum_{s=0}^{s_e} T_s \epsilon_s + \frac{T_s}{T^2}\right) \quad \text{By Equation (40)}$$

$$\leq \tilde{O}_T(T\nu_T).$$

$\square$

## D.3   Proof of Proposition 5 (Concentration of Cost)

*proof.* The following lemma is a result of McDiarmid's inequality and shows that the random variable corresponding to $T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W)$ concentrates around a conditional expectation.

**Lemma 6.** *Under Assumptions 1–8 , for every $s \in [0 : s_e]$ there exists an event $E_s^{\mathrm{M}}$ such that $E_s^{\mathrm{M}}$ depends only on the random variables in $W_s$ and $\hat{\theta}_s$, such that $E_s^{\mathrm{M}} \subseteq \{\forall t \in [T_s : T_{s+1} - 1], |w_t| \leq \log^2(T)\}$, and such that conditional on $E_2^s$, $\mathbb{P}(E_s^{\mathrm{M}} \mid \hat{\theta}_s) \geq 1 - o_T(1/T^8)$ and for $\epsilon \geq 1/T$ and for sufficiently large $T$,*

$$\mathbb{P}\left(\left|T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) - \mathbb{E}\left[T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \,\Big|\, E_s^{\mathrm{M}}, \hat{\theta}_s\right]\right| \geq \epsilon \,\Big|\, \hat{\theta}_s\right)$$

$$\leq \frac{1}{T^8} + 2\exp\left(-\frac{\epsilon^2}{2 T_s c^2}\right)$$

*for some $c = \tilde{O}_T(1)$.*

The proof of Lemma 6 can be found in Appendix E.4. We also want that taking expectation conditional on $E_s^{\mathrm{M}}$ does not significantly change the expected cost.

**Lemma 7.** *Under Assumptions 1–8, if $E_s^{\mathrm{M}} \subseteq \{\forall t \in [T_s : T_{s+1} - 1], |w_t| \leq \log^2(T)\}$ and conditional on event $E_2^s$ we have $\mathbb{P}(E_s^{\mathrm{M}}) \geq 1 - o_T(1/T^8)$, then conditional on event $E_2^s$,*

$$\mathbb{E}\left[T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \,\Big|\, \hat{\theta}_s\right] \overset{a.s.}{\geq} \mathbb{E}\left[T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \,\Big|\, E_s^{\mathrm{M}}, \hat{\theta}_s\right] - \tilde{O}_T(1),$$

(41)

*where the term $\tilde{O}_T(1)$ does not depend on s.*

The proof of Lemma 7 can be found in Appendix E.5. Combining Lemma 6 for $\epsilon = c\sqrt{T_s}\log(T)$ and Lemma 7 for sufficiently large $T$, we have the following conditional on event $E_2^s$:

$$\mathbb{P}\left(T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) - \mathbb{E}\left[T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \,\Big|\, \hat{\theta}_s\right] \geq c\sqrt{T_s}\log(T) + \tilde{O}_T(1) \,\Big|\, \hat{\theta}_s\right)$$
$$\leq \frac{1}{T^8} + 2\exp\left(-\frac{\log^2(T)}{2}\right). \tag{42}$$

Now applying a union bound over all $s \in [0 : s_e]$ gives the following result:

$$\mathbb{P}\left(\sum_{s=0}^{s_e} T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) - \sum_{s=0}^{s_e}\mathbb{E}\left[T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \,\Big|\, \hat{\theta}_s\right] \geq \sum_{s=0}^{s_e}\left(c\sqrt{T_s}\log(T) + \tilde{O}_T(1)\right)\right)$$
$$\leq \mathbb{P}\left(\exists s \in [0 : s_e] : T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) - \mathbb{E}\left[T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \,\Big|\, \hat{\theta}_s\right] \geq c\sqrt{T_s}\log(T) + \tilde{O}_T(1)\right)$$
$$\leq \sum_{s=0}^{s_e}\mathbb{P}\left(T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) - \mathbb{E}\left[T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \,\Big|\, \hat{\theta}_s\right] \geq c\sqrt{T_s}\log(T) + \tilde{O}_T(1)\right)$$
$$\leq \sum_{s=0}^{s_e}\left(\frac{1}{T^8} + 2\exp\left(-\frac{\log^2(T)}{2}\right) + \mathbb{P}(\neg E_2^s)\right) \qquad\qquad \text{Equation (42)}$$
$$\leq \tilde{O}_T\left(\frac{1}{T^2}\right). \qquad\qquad\qquad\qquad \text{Equation (39)} \tag{43}$$

Note that

$$\sum_{s=0}^{s_e} c\sqrt{T_s}\log(T) = \tilde{O}_T(\sqrt{T}), \tag{44}$$

therefore combining Equations (44) and (43), we have that

$$\mathbb{P}\left(\sum_{s=0}^{s_e} T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) - \sum_{s=0}^{s_e}\mathbb{E}\left[T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \,\Big|\, \hat{\theta}_s\right] \geq \tilde{O}_T(\sqrt{T})\right)$$
$$\leq \tilde{O}_T\left(\frac{1}{T^2}\right). \tag{45}$$

Equation (45) differs from the desired result of Proposition 5 in that the first summation is over trajectories starting at position 0 as opposed to $\hat{x}_{T_s}$. Therefore, the last part of this proof is to bound

$$\left|\sum_{s=0}^{s_e} T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, \hat{x}_{T_s}, W_s) - \sum_{s=0}^{s_e} T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s)\right|.$$

To do this, we will use the following lemma that is a consequence of Assumption 8.

**Lemma 8.** *Under Assumptions 1–6 and 8, if $\|\theta - \theta^*\|_\infty = \epsilon \leq \epsilon_{A8}$, then for any $K \in (K_L^\theta, K_U^\theta)$, $t \leq T$, and $|x|, |y| \leq 4\log^2(T)$ and any noise random variables $W'$, conditional on event $E_{A8}(C_K^\theta, W')$,*

$$\left|t \cdot J(\theta^*, C_K^\theta, t, x, W') - t \cdot J(\theta^*, C_K^\theta, t, y, W')\right| = \tilde{O}_T\left(|x - y| + \epsilon\right).$$

The proof of Lemma 8 can be found in Appendix E.6.

In order to use Lemma 8, we must show that $|\hat{x}_{T_s}| \le 4\log^2(T)$. Recall that $\hat{x}_{T_s}$ is the position at time $T_s$ if the position at time $T_0$ is $\hat{x}_{T_0} = x'_{T_0}$, where $x'_{T_0}$ is the position of the controller $C^{\mathrm{alg}}$ at time $T_0$. Because $E_{\mathrm{safe}} \subseteq E$, under event $E$ we have that $C^{\mathrm{alg}}$ is safe for dynamics $\theta^*$. Therefore by Lemma 4, $|x'_{T_0}| \le 4\log^2(T)$. Because $E_2 \subseteq E$, under event $E$ we also have that $\|\hat{\theta}_s - \theta^*\|_\infty \le \tilde{O}_T(\nu_T)$ for all $s \in [0:s_e]$ and sufficiently large $T$. Therefore, since $\hat{x}_{T_0} = x'_{T_0}$ and the control $C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}(x)$ is safe with respect to $\hat{\theta}_s$ for any $x$, again by Lemma 4 we have that under event $E$ and for sufficiently large $T$, $|\hat{x}_{T_s}| \le 4\log^2(T)$. Now we can apply Lemma 8 to get that, conditional on event $E \cap \bigcap_{s=0}^{s_e} E_{\mathrm{A8}}(C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, W_s)$,

$$
\left| \sum_{s=0}^{s_e} T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, \hat{x}_{T_s}, W_s) - \sum_{s=0}^{s_e} T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \right|
$$

$$
\le \sum_{s=0}^{s_e} \left| T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, \hat{x}_{T_s}, W_s) - T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \right|
$$

$$
\le \sum_{s=0}^{s_e} \tilde{O}_T \left( \hat{x}_{T_s} + \|\hat{\theta}_s - \theta^*\|_\infty \right)
$$

$$
\le \tilde{O}_T(1). \tag{46}
$$

A union bound gives that $\mathbb{P}(\bigcap_{s=0}^{s_e} E_{\mathrm{A8}}(C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, W_s)) = 1 - o_T(1/T^2)$. Combining Equation (45) with Equation (46) with a union bound gives that conditional on event $E$ with probability $1 - o_T(1/T)$,

$$
\sum_{s=0}^{s_e} T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, \hat{x}_{T_s}, W_s) - \sum_{s=0}^{s_e} \mathbb{E}\left[ T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \,\Big|\, \hat{\theta}_s \right] \le \tilde{O}_T(\sqrt{T}),
\tag{47}
$$

which is the desired result of Proposition 5. $\square$

## D.4 Proof of Proposition 6 (Regret of Enforcing Safety)

*proof.* Intuitively, $R_3$ is the regret caused by enforcing safety and deviating from the controller $C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}$. Lemma 9 bounds the cost of deviating from $C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}$ as a sum over all times the algorithm deviates.

**Lemma 9.** *Recall $u_t^{\mathrm{safeU}}$ and $u_t^{\mathrm{safeL}}$ defined in Algorithm 2 Lines 11 and 12. Let $X_t^U$ and $X_t^L$ be the indicators for the events that at time $t$, $C^{\mathrm{alg}}(x'_t) = u_t^{\mathrm{safeU}}$ or $C^{\mathrm{alg}}(x'_t) = u_t^{\mathrm{safeL}}$, respectively. Under Assumptions 1–8 and conditional on event $E$, with probability $1 - o_T(1/T)$*

$$
\sum_{s=0}^{s_e} T_s J(\theta^*, C^{\mathrm{alg}}, T_s, x'_{T_s}, W_s) - \sum_{s=0}^{s_e} T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, \hat{x}_{T_s}, W_s)
$$

$$
\le \tilde{O}_T \left( \sum_{s=0}^{s_e} \epsilon_s T_s \right) + \sum_{s=0}^{s_e} \sum_{t=T_s}^{T_{s+1}-1} X_t^U \cdot \tilde{O}_T \left( \left| u_t^{\mathrm{safeU}} - C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}(x'_t) \right| \right) + X_t^L \cdot \tilde{O}_T \left( \left| u_t^{\mathrm{safeL}} - C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}(x'_t) \right| \right).
$$

The proof of Lemma 9 can be found in Appendix E.7. We also remind the reader that the coefficients of the $\tilde{O}_T(\cdot)$ terms in Lemma 9 do not depend on $t$ or $s$, and are a function of known problem parameters and $\log(T)$ factors. The next tool we need is to be able to bound the difference in control when applying safety in Algorithm 2 compared to the control when not applying safety. We can do that as follows.

**Lemma 10.** *Under Assumptions 1–8 and conditional on event $E$, for any $t$ such that $1/\nu_T^2 \leq t \leq T$, if $s = \lfloor \log_2(t\nu_T^2) \rfloor$ and $u_t^{\mathrm{safeU}} \leq C_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}(x_t')$ (which is equivalent to $C^{\mathrm{alg}}(x_t') = u_t^{\mathrm{safeU}}$), then,*

$$|u_t^{\mathrm{safeU}} - C_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}(x_t')| \leq \tilde{O}_T(\epsilon_s). \tag{48}$$

*Similarly, if $u_t^{\mathrm{safeL}} \geq C_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}(x_t')$, then conditional on event $E$,*

$$|u_t^{\mathrm{safeL}} - C_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}(x_t')| \leq \tilde{O}_T(\epsilon_s). \tag{49}$$

The proof of Lemma 10 can be found in section E.8. Combining Lemmas 9 and 10, we have that conditional on event $E$, with probability $1 - o_T(1/T)$,

$$
\begin{aligned}
R_3 &= \sum_{s=0}^{s_e} T_s J(\theta^*, C^{\mathrm{alg}}, T_s, x_{T_s}', W_s) - \sum_{s=0}^{s_e} T_s J(\theta^*, C_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}, T_s, \hat{x}_{T_s}, W_s) \\
&\leq \tilde{O}_T\left(\sum_{s=0}^{s_e} \epsilon_s T_s\right) + \sum_{s=0}^{s_e} \sum_{t=T_s}^{T_{s+1}-1} \left(X_t^U \cdot \tilde{O}_T\left(\left|u_t^{\mathrm{safeU}} - C_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}(x_t')\right|\right)\right. \\
&\qquad \left. + X_t^L \cdot \tilde{O}_T\left(\left|u_t^{\mathrm{safeL}} - C_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}(x_t')\right|\right)\right) \qquad\qquad \text{Lemma 9} \\
&\leq \tilde{O}_T\left(\sum_{s=0}^{s_e} \epsilon_s T_s\right) + \sum_{s=0}^{s_e} \sum_{t=T_s}^{T_{s+1}-1} X_t^U \cdot \tilde{O}_T(\epsilon_s) + X_t^L \cdot \tilde{O}_T(\epsilon_s) \qquad \text{Lemma 10} \\
&\leq \tilde{O}_T(T\nu_T) + \sum_{t=1/\nu_T^2}^{T-1} X_t^U \cdot \tilde{O}_T(\nu_T) + X_t^L \cdot \tilde{O}_T(\nu_T) \qquad\qquad E_2 \subseteq E \\
&\leq \tilde{O}_T(T\nu_T) + \tilde{O}_T(T\nu_T) \\
&= \tilde{O}_T(T\nu_T).
\end{aligned}
$$

The key application of event $E$ in the above result is that $E_2 \subseteq E$ implies that under event $E$, $\max_{s \in [0:s_e]} \epsilon_s = \tilde{O}_T(\nu_T)$. $\qquad\square$

# E Proofs of Lemmas from Appendix D

## E.1 Proof of Lemma 3

Recall the notation that $x_t'$ is the position at time $t$ when using the controller $C^{\mathrm{alg}}$. We will prove the following. For sufficiently large $T$ and any $t \in [T_0 : T]$, if $C^{\mathrm{alg}}(x_{t-1}')$ is safe for

dynamics $\theta^*$, then conditional on $E_1 \cap E_2$, we have that both $u_t^{\text{safeL}} \leq u_t^{\text{safeU}}$ and $C^{\text{alg}}(x_t')$ is safe for dynamics $\theta^*$. Because we assume in this lemma that $u_{T_0-1} = C^{\text{alg}}(x_{T_0-1}')$ is safe with respect to dynamics $\theta^*$, this will prove by induction the desired result that conditional on $E_1 \cap E_2$ and for sufficiently large $T$, $u_t^{\text{safeL}} \leq u_t^{\text{safeU}}$ for all $t \in [T_0 : T]$.

Fix a given $t$, and define $s = \lfloor \log_2(t/\nu_T^2) \rfloor$. Assume $C^{\text{alg}}(x_{-1}')$ is safe for dynamics $\theta^*$. Then under event $E_1$, we have that $|x_t'| \leq \|D\|_\infty + |w_{t-1}| \leq B_x$. Let $v = \frac{D_U - a^* x_t' - 4\epsilon_s B_x}{b^*}$. We will show that $u_t^{\text{safeU}} \geq v$. Note that $a^* x_t' + b^* v = D_U - 4\epsilon_s B_x$. For sufficiently large $T$, because $D_U - D_L \geq \frac{1}{\log(T)}$ (Assumption 3) and $\epsilon_s = \tilde{O}_T(\nu_T) = o_T(1/\log(T))$ under $E_1 \cap E_2$, this implies that

$$D_L \leq a^* x_t' + b^* v \leq D_U.$$

Therefore $v$ is safe for dynamics $\theta^*$, which implies by Lemma 4 that under event $E_1$ and for sufficiently large $T$,

$$|v| \leq B_x.$$

Under event $E_1 \cap E_2$, $\|\theta^* - \hat{\theta}_s\|_\infty \leq \epsilon_s$, therefore by the above results we have that under $E_1 \cap E_2$ and for sufficiently large $T$,

$$\max_{\|\hat{\theta}_s - \theta\|_\infty \leq \epsilon_s} a x_t' + b v \leq a^* x_t' + b^* v + 2\epsilon_s |x_t'| + 2\epsilon_s |v|$$

$$\leq a^* x_t' + b^* v + 4\epsilon_s B_x \qquad\qquad |v| \leq B_x, |x_t'| \leq B_x$$

$$= D_U. \qquad\qquad\qquad\qquad \text{Def of } v$$

This implies by the definition of $u_t^{\text{safeU}}$ that

$$u_t^{\text{safeU}} \geq v = \frac{D_U - a^* x_t' - 4\epsilon_s B_x}{b^*}.$$

By the same logic, we also have that

$$u_t^{\text{safeL}} \leq \frac{D_L - a^* x_t' + 4\epsilon_s B_x}{b^*}.$$

For sufficiently large $T$ under event $E_2$, $\frac{8\epsilon_s B_x}{b^*} = \tilde{O}_T(\nu_T) \leq \frac{1}{\log(T)}$. Therefore, using that $D_U \geq D_L + \frac{1}{\log(T)}$ by Assumption 3, we can conclude that under event $E_1 \cap E_2$ for sufficiently large $T$,

$$u_t^{\text{safeL}} \leq \frac{D_L - a^* x_t' + 4\epsilon_s B_x}{b^*}$$

$$\leq \frac{D_U - a^* x_t' - 4\epsilon_s B_x}{b^*}$$

$$\leq u_t^{\text{safeU}}.$$

This implies that $u_t^{\text{safeL}} \leq C^{\text{alg}}(x_t') \leq u_t^{\text{safeU}}$, which by construction under event $E_1 \cap E_2$ implies that $D_L \leq a^* x_t' + b^* C^{\text{alg}}(x_t') \leq D_U$. Finally, this gives that $C^{\text{alg}}(x_t')$ is safe for dynamics $\theta^*$. Therefore, we have shown the two desired results that $u_t^{\text{safeL}} \leq u_t^{\text{safeU}}$ and $C^{\text{alg}}(x_t')$ is safe for dynamics $\theta^*$.

As mentioned above, this implies by induction the desired result that $u_t^{\text{safeL}} \leq u_t^{\text{safeU}}$ for all $t \in [T_0, T]$ conditional on $E_1 \cap E_2$ as long as $C^{\text{alg}}(x_{T_0-1}')$ is safe with respect to $\theta^*$.

## E.2 Proof of Lemma 4 (Bounded positions and controls)

*proof.* Define $\gamma_T = \max_{t \in [T]} \|\theta^* - \theta_t\|_\infty$, and we know that $\gamma_T \le \frac{1}{\log(T)}$ by assumption. At time $t$, the control used by controller $C_t$ is safe for dynamics $\theta_t$ by assumption of the lemma, so by Definition 2, for all $t$, if $u_t = C_t(x_t)$ then

$$D_{\mathrm{L}} \le a_t x_t + b_t u_t \le D_{\mathrm{U}}. \tag{50}$$

By definition of $\gamma_T$, this implies that

$$D_{\mathrm{L}} - \gamma_T |x_t| - \gamma_T |u_t| \le a^* x_t + b^* u_t \le D_{\mathrm{U}} + \gamma_T |x_t| + \gamma_T |u_t|. \tag{51}$$

The right inequality in Equation (51) implies that

$$b^* u_t - \gamma_T |u_t| \le D_{\mathrm{U}} + \gamma_T |x_t| - a^* x_t,$$

which for $u_t \ge 0$ implies that $|u_t| \le \frac{\|D\|_\infty + a^* |x_t| + \gamma_T |x_t|}{b^* - \gamma_T}$. The left inequality in Equation (51) implies the same for $u_t \le 0$, and therefore we have that Equation (51) implies that

$$|u_t| \le \frac{\|D\|_\infty + a^* |x_t| + \gamma_T |x_t|}{b^* - \gamma_T}. \tag{52}$$

First we prove Equations (33) and (34) by induction.

**Base Case:** At time $t = 0$, we have by assumption that $|x_0| \le 4 \log^2(T)$. Furthermore, Equation (52) implies that

$$
\begin{aligned}
|u_0| &\le \frac{\|D\|_\infty + a^* |x_0| + \gamma_T |x_0|}{b^* - \gamma_T} && \text{Equation (52)} \\
&\le \frac{\|D\|_\infty + (a^* + \gamma_T) 4 \log^2(T)}{b^* - \frac{1}{\log(T)}} && \text{Equation (32)} \\
&\le \frac{\log^2(T) + (a^* + \frac{1}{\log(T)}) 4 \log^2(T)}{b^* - \frac{1}{\log(T)}} && \text{Assumption 3} \\
&\le \frac{\log^2(T) + (a^* + b^*/2) 4 \log^2(T)}{b^*/2} && \text{Sufficiently large } T \\
&\le \frac{2(1 + 4a^* + 2b^*) \log^2(T)}{b^*} && (53) \\
&< B_x,
\end{aligned}
$$

for $T$ sufficiently large such that $2(1 + 4a^* + 2b^*)/b^* \le \log(T)$ and $1/\log(T) \le b^*/2$.

**Induction Hypothesis:** Assume Equations (33) and (34) are true for all times less than or equal to $t$.

**Induction Step:** Now we will prove that Equations (33) and (34) hold at time $t+1$.

$$|x_{t+1}| = |a^* x_t + b^* u_t + w_t|$$
$$= |a_t x_t + b_t u_t + w_t + (a^* - a_t)x_t + (b^* - b_t)u_t|$$
$$\leq |a_t x_t + b_t u_t| + |w_t| + |(a^* - a_t)x_t| + |(b^* - b_t)u_t| \quad \text{Triangle Inequality}$$
$$\overset{\text{a.s.}}{\leq} \|D\|_\infty + \log^2(T) + \gamma_T |x_t| + \gamma_T |u_t| \quad \text{Equation (32), Equation (50), event } E_1$$
$$\leq \|D\|_\infty + \frac{1}{\log(T)}(|x_t| + |u_t|) + \log^2(T) \quad \text{Equation (32)}$$
$$\leq \|D\|_\infty + \frac{2}{\log(T)} B_x + \log^2(T) \quad \text{Ind. Hyp.}$$
$$\leq \|D\|_\infty + 3\log^2(T)$$
$$\leq 4\log^2(T) \quad \text{Assumption 3}$$
$$< \log^3(T)$$
$$= B_x.$$

Above we need $T$ large enough such that $\log(T) > 4$. Since we showed that $|x_{t+1}| \leq 4\log^2(T)$, this also implies by Equations (52) and (53) that for sufficiently large $T$, $|u_{t+1}| < B_x$. Therefore we have shown Equations (33) and (34) for time $t+1$, completing the induction proof.

Now we will prove Equations (35) and (36) with a similar proof by induction. If the controller $C_t$ is non-random and $x_0$ is not random, this implies that $\mathbb{E}[|x_0|] = |x_0| \leq 4\log^2(T)$ and $\mathbb{E}[|u_0|] = |u_0| \leq \frac{2(5+4a^*)\log^2(T)}{b^*}$ by Equation (53). This proves the base case. For the inductive step, we have that

$$\mathbb{E}[|x_{t+1}|]$$
$$= \mathbb{E}[|a^* x_t + b^* u_t + w_t|]$$
$$\leq \mathbb{E}[|a_t x_t + b_t u_t|] + \mathbb{E}[|w_t|] + \mathbb{E}[|(a^* - a_t)x_t|] + \mathbb{E}[|(b^* - b_t)u_t|] \quad \text{Triangle Inequality}$$
$$\leq \|D\|_\infty + \log^2(T) + \gamma_T \mathbb{E}[|x_t|] + \gamma_T \mathbb{E}[|u_t|] \quad \text{Equations (32), (50), } w_t \text{ sub-Gaussian}$$
$$\leq \|D\|_\infty + \frac{1}{\log(T)}(\mathbb{E}[|x_t|] + \mathbb{E}[|u_t|]) + \log^2(T) \quad \text{Equation (32)}$$
$$\leq \|D\|_\infty + \frac{2}{\log(T)} B_x + \log^2(T) \quad \text{Ind. Hyp.}$$
$$\leq \|D\|_\infty + 3\log^2(T)$$
$$\leq 4\log^2(T) \quad \text{Assumption 3}$$
$$< \log^3(T)$$
$$= B_x.$$

We have shown that $\mathbb{E}[|x_{t+1}|] \leq 4\log^2(T)$, therefore by Equation (52) and the same algebraic

steps as used in Equation (53), we have that for sufficiently large $T$,

$$\mathbb{E}[|u_{t+1}|] \leq \frac{\|D\|_\infty + a^* \, \mathbb{E}[|x_{t+1}|] + \gamma_T \, \mathbb{E}[|x_{t+1}|]}{b^* - \gamma_T}$$

$$\leq \frac{\|D\|_\infty + (a^* + \gamma_T)4\log^2(T)}{b^* - \frac{1}{\log(T)}}$$

$$\leq \frac{2(1 + 4a^* + 2b^*)\log^2(T)}{b^*}$$

$$< B_x.$$

This completes the second proof by induction, proving Equations (35) and (36). $\qquad\square$

## E.3   Proof of Lemma 5

*proof.* For this proof, we need the following version of Lemma 8 that applies for expectations rather than with high probability.

**Lemma 11.** *Let $x, y$ be two random variables independent of noises $W' = \{w'_i\}_{i=0}^{t-1}$ such that for some $L = \tilde{O}_T(1)$, both $\mathbb{P}(|x| \geq L)\,\mathbb{E}[x^2 \mid |x| \geq L] = o_T\left(\frac{1}{T^{10}}\right)$ and $\mathbb{P}(|y| \geq L)\,\mathbb{E}[y^2 \mid |y| \geq L] = o_T\left(\frac{1}{T^{10}}\right)$ and $\mathbb{P}(|x| \leq 4\log^2(T)) = 1 - o_T(1/T^{11})$ and $\mathbb{P}(|y| \leq 4\log^2(T)) = 1 - o_T(1/T^{11})$. Then under Assumptions 1–6 and 8, if $\|\theta - \theta^*\|_\infty = \epsilon \leq \epsilon_{A8}$, then for any $K \in (K_L^\theta, K_U^\theta)$ and $t \leq T$,*

$$\left| \mathbb{E}\left[ t \cdot J(\theta^*, C_K^\theta, t, x, W') - t \cdot J(\theta^*, C_K^\theta, t, y, W') \right] \right| = \tilde{O}_T\left( \mathbb{E}[|x - y|] + \epsilon + \frac{1}{T^2} \right). \tag{54}$$

The proof of Lemma 11 can be found in Appendix E.10. We also need the following generalization of Lemma 4, which bounds the positions for any starting position $x$.

**Lemma 12.** *Let $x_0, x_1, ...x_T$ be the sequences of positions when starting at position $x_0 = x$ and using controller $C_t$ at time $t$. Suppose that the control $C_t(x_t)$ is safe for dynamics $\theta_t$ and $\|\theta_t - \theta^*\| \leq \frac{1}{\log(T)}$ for all $t < T$. For sufficiently large $T$ under Assumption 3,*

$$\forall t \leq T, \ |x_t| = O_T(|x| + \|D\|_\infty + \max_{i \leq t-1} |w_i|).$$

$$\forall t < T, \ |C_t(x_t)| = O_T(|x| + \|D\|_\infty + \max_{i \leq t-1} |w_i|).$$

The proof of Lemma 12 can be found in Appendix E.11.

Because $C_{K^*}^{\theta^*}, \{C_{K_s^*}^{\theta^*}\}_{s=0}^{s_e}$ are safe for dynamics $\theta^*$, the sequence $x_0^*, x_1^*, ...$ starts at $x_0^* = 0$, and $\|D\|_\infty \leq \log^2(T)$ by Assumption 3, Lemma 12 implies that

$$|x_{T_s}^*| = O_T\left( \max_{i \leq T_s - 1} |w_i| + \log^2(T) \right). \tag{55}$$

**Lemma 13.** *Suppose $w_t$ for $t < T$ are sub-Gaussian and $F$ is an event such that $\mathbb{P}(F) = 1 - o_T(1/T^{11})$. Then*

$$\mathbb{E}[\max_{i \leq t} w_i^2 \mid \neg F]\mathbb{P}(\neg F) = o_T\left( \frac{1}{T^{10}} \right).$$

The proof of Lemma 13 can be found in Appendix E.12. Define $F = \{|x^*_{T_s}| < \log^3(T)\}$. Event $E_1$ implies $F$ by Lemma 4, and therefore $\mathbb{P}(F) \geq \mathbb{P}(E_1) = 1 - o_T(1/T^{11})$. Therefore, we have by Equation (55) that

$$
\begin{aligned}
&\mathbb{P}(\neg F)\,\mathbb{E}[|x^*_{T_s}|^2 \mid \neg F] \\
&= O_T\left(\mathbb{P}(\neg F)\,\mathbb{E}\left[\max_{i \leq T_s-1} w_i^2 \,\Big|\, \neg F\right]\right) + \tilde{O}_T\left(\mathbb{P}(\neg F)\right) \quad \text{[Eq. (55) and } (a+b)^2 \leq 2a^2 + 2b^2] \\
&= o_T\left(\frac{1}{T^{10}}\right). \hspace{5.5cm} \text{Lemma 13, } \mathbb{P}(\neg F) = o_T(1/T^{11}) \quad (56)
\end{aligned}
$$

Also, note that Lemma 4 implies that $\mathbb{P}(x^*_{T_s} \leq 4\log^2(T)) \geq \mathbb{P}(E_1) = 1 - o_T(1/T^{11})$. We can therefore apply Lemma 11 with $x = x^*_{T_s}, y = 0, L = \log^3(T), \epsilon = 0$. Applying Lemma 11 gives the following desired result.

$$
\begin{aligned}
&\mathbb{E}\left[\left|T_s J(\theta^*, C^{\theta^*}_{K^*_s}, T_s, x^*_{T_s}, W_s) - T_s J(\theta^*, C^{\theta^*}_{K^*_s}, T_s, 0, W_s)\right|\right] \\
&= \tilde{O}_T\left(\mathbb{E}\left[|x^*_{T_s}|\right] + \frac{1}{T^2}\right) \hspace{4cm} \text{Lemma 11} \\
&= \tilde{O}_T(1). \hspace{6cm} \text{Lemma 4 for sufficiently large } T
\end{aligned}
$$

Note that we can apply the expectation form of Lemma 4 in the second inequality above because $(C^{\theta^*}_{K^*}, \{C^{\theta^*}_{K^*_s}\}^{s_e}_{s=0})$ are non-random controllers. $\hspace{3cm}\square$

## E.4 Proof of Lemma 6 (Concentration of Conditional Expected Cost)

*proof.* We will use the following form of McDiarmid's Inequality for high probability events.

**Lemma 14** (McDiarmid's Inequality Combes [2015]). *Let $f$ be a function such that $f : \mathcal{X}_1 \times \mathcal{X}_2... \times \mathcal{X}_n \to \mathbb{R}$ and let $\mathcal{Y} \in \mathcal{X}_1 \times \mathcal{X}_2... \times \mathcal{X}_n$ be a subset of the domain such that for some $c$, if $(x_1, ..., x_n), (x'_1, ..., x'_n) \in \mathcal{Y}$, then*

$$
|f(x_1, ..., x_n) - f(x'_1, ..., x'_n)| \leq \sum_{i:x_i \neq x'_i} c.
$$

*Let $X_1, X_2, ..., X_n$ be independent random variables and $X_i \in \mathcal{X}_i$ for all $i$. Define $p = 1 - \mathbb{P}((X_1, ..., X_n) \in \mathcal{Y})$ and let $m = \mathbb{E}[f(X_1, ..., X_n) \mid (X_1, ..., X_n) \in \mathcal{Y}]$. Then for any $\epsilon > 0$,*

$$
\mathbb{P}(|f(X_1, ..., X_n) - m| \geq \epsilon) \leq 2p + 2\exp\left(-\frac{2\max(0, \epsilon - pnc)^2}{nc^2}\right).
$$

Define the function $f_{\hat{\theta}_s}(W_m)$ as

$$
f_{\hat{\theta}_s}(W_m) = T_s J\left(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_m\right).
$$

We want to apply McDiarmid's Inequality to $f_{\hat{\theta}_s}$ conditional on $\hat{\theta}_s$ when $E_2^s$ holds, which requires the following bounded difference result.

**Lemma 15.** *Under Assumptions 1–8, given $\hat{\theta}_s$ there exists a fixed $\mathcal{Y}_s \in [-\log^2(T), \log^2(T)]^{T_s}$ such that the event $E_s^M := \{W_s \in \mathcal{Y}_s\}$ satisfies $\mathbb{P}(E_s^M \mid \hat{\theta}_s) \geq 1 - o_T(1/T^8)$, and conditional on $\hat{\theta}_s$ and $E_2^s$, if $E_s^M$ holds when $W_s = \{w_i\}_{i=T_s}^{T_{s+1}-1}$ and when $W_s' = \{w_i'\}_{i=T_s}^{T_{s+1}-1}$, then*

$$\left| f_{\hat{\theta}_s}(W_s) - f_{\hat{\theta}_s}(W_s') \right| \leq \sum_{i=T_s, w_i \neq w_i'}^{T_{s+1}-1} c$$

*for some $c = \tilde{O}_T(1)$.*

The proof of Lemma 15 can be found in Appendix E.9. We will now apply Lemma 14 for the function $f_{\hat{\theta}_s}$ conditional on $\hat{\theta}_s$ and $E_2^s$ using Lemma 15. Conditional on $E_2^s$ (where $c$ is from Lemma 15), the following holds for $\epsilon \geq 1/T$ and $T$ sufficiently large.

$$\mathbb{P}\left( \left| f_{\hat{\theta}_s}(W_s) - \mathbb{E}[f_{\hat{\theta}_s}(W_s) \mid E_s^M] \right| \geq \epsilon \;\Big|\; \hat{\theta}_s \right)$$

$$\leq 2\mathbb{P}(\neg E_s^M \mid \hat{\theta}_s) + 2\exp\left( -\frac{2\max\left(0, \epsilon - cT_s\mathbb{P}(\neg E_s^M \mid \hat{\theta}_s)\right)^2}{T_s c^2} \right)$$

$$= o_T\left(\frac{1}{T^8}\right) + 2\exp\left(-\frac{\epsilon^2}{2T_s c^2}\right) \qquad [\epsilon \geq 1/T,\ \mathbb{P}(\neg E_s^M | \hat{\theta}_s) = o_T(1/T^8),\ \text{suff. large } T\ ]$$

$$\leq \frac{1}{T^8} + 2\exp\left(-\frac{\epsilon^2}{2T_s c^2}\right). \qquad [\text{Suff. large } T]$$

$\square$

## E.5   Proof of Lemma 7 (Unconditional Cost vs Conditional Cost)

*proof.* By the Law of Total Expectation,

$$\mathbb{E}\left[ T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \;\Big|\; \hat{\theta}_s \right]$$

$$= \mathbb{E}\left[ T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \;\Big|\; E_s^M, \hat{\theta}_s \right] \mathbb{P}(E_s^M \mid \hat{\theta}_s)$$

$$\qquad + \mathbb{E}\left[ T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \;\Big|\; \neg E_s^M, \hat{\theta}_s \right] \mathbb{P}(\neg E_s^M \mid \hat{\theta}_s)$$

$$\geq \mathbb{E}\left[ T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \;\Big|\; E_s^M, \hat{\theta}_s \right] \mathbb{P}(E_s^M | \hat{\theta}_s) \qquad \text{Cost is non-negative}$$

$$= \mathbb{E}\left[ T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \;\Big|\; E_s^M, \hat{\theta}_s \right] - \mathbb{E}\left[ T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \;\Big|\; E_s^M, \hat{\theta}_s \right] \mathbb{P}(\neg E_s^M | \hat{\theta}_s)$$

$$= \mathbb{E}\left[ T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \;\Big|\; E_s^M, \hat{\theta}_s \right] - o_T\left(\frac{1}{T}\right) \mathbb{E}\left[ T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \;\Big|\; E_s^M, \hat{\theta}_s \right]$$

$$= \mathbb{E}\left[ T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \;\Big|\; E_s^M, \hat{\theta}_s \right] - o_T((q+r)B_x^2)$$

$$= \mathbb{E}\left[ T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \;\Big|\; E_s^M, \hat{\theta}_s \right] - \tilde{O}_T(1).$$

To see the step from the 5th to the 6th line, note that $E_s^M \subseteq \{\forall t \in [T_s : T_{s+1} - 1], |w_t| \leq \log^2(T)\}$ by assumption and that $E_2^s$ implies that for sufficiently large $T$, $\|\theta^* - \hat{\theta}_s\| \leq \frac{1}{\log(T)}$, therefore by Lemma 4 we have that the magnitudes of the positions and controls are all

44

bounded by $B_x$ conditional on events $E_2^s$ and $E_s^M$. Therefore, the cost at each time step conditional on these events is at most $(q+r)B_x^2$, which gives that conditional on event $E_2^s$,

$$\mathbb{E}\left[T_s J(\theta^*, C_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}, T_s, 0, W_s) \mid E_s^M, \hat{\theta}_s\right]$$
$$\leq T_s(q+r)B_x^2 \qquad E_2^s, E_s^M \subseteq \{\forall t \in [T_s : T_{s+1}-1], |w_t| \leq \log^2(T)\}, \text{Lemma 4}$$
$$\leq T(q+r)B_x^2.$$

$\square$

## E.6  Proof of Lemma 8

*proof.* If $|x - y| \leq \delta_{A8}$ then this follows directly from Assumption 8. Now for the rest of this proof assume $|x - y| > \delta_{A8}$ and WLOG assume $x \leq y$. Choose $\delta$ to be the largest real number satisfying $\delta \leq \delta_{A8}$ such that $\frac{|x-y|}{\delta}$ is an integer. Because $\delta_{A8} < |x - y|$, there must exist an integer in the range $\left[\frac{|x-y|}{\delta_{A8}}, \frac{2|x-y|}{\delta_{A8}}\right]$. Therefore, $\delta \geq \delta_{A8}/2 = \tilde{\Omega}_T(1)$ by definition of $\delta_{A8}$. Because $|x|, |y| < 4\log^2(T)$ and $x \leq y$, we know that for all $i \in [0 : \frac{|x-y|}{\delta}]$, we have $|x + i\delta| \leq 4\log^2(T)$. For $i \in [0 : \frac{|x-y|}{\delta} - 1]$, by Assumption 8, under event $E_{A8}(C_K^\theta, W')$

$$|t \cdot J(\theta^*, C_K^\theta, t, x + i\delta, W') - t \cdot J(\theta^*, C_K^\theta, t, x + (i+1)\delta, W')| = \tilde{O}_T(\delta + \epsilon).$$

By the triangle inequality, this implies that conditional on event $E_{A8}(C_K^\theta, W')$,

$$\left|t \cdot J(\theta^*, C_K^\theta, t, x, W') - t \cdot J(\theta^*, C_K^\theta, t, y, W')\right|$$
$$\leq \sum_{i=0}^{\frac{|x-y|}{\delta}-1} \left|t \cdot J(\theta^*, C_K^\theta, t, x + i\delta, W') - t \cdot J(\theta^*, C_K^\theta, t, x + (i+1)\delta, W')\right|$$
$$= \tilde{O}_T\left(\frac{|x-y|}{\delta}(\delta + \epsilon)\right)$$
$$= \tilde{O}_T\left(|x-y| + \frac{8\log^2(T)}{\delta}\epsilon\right) \qquad\qquad |x|, |y| < 4\log^2(T)$$
$$= \tilde{O}_T\left(|x-y| + \epsilon\right). \qquad\qquad \delta = \tilde{\Omega}_T(1)$$

$\square$

## E.7  Proof of Lemma 9 (Cost of safety controls)

*proof.* The first tool for this proof is the following lemma, which informally states that being off by a small amount of control has a small impact on the overall cost.

**Lemma 16.** *Under Assumptions 1–8, with conditional probability $1 - o_T(1/T)$ given event $E$, for all $s \in [0 : s_e]$,*

$$|T_s \cdot J(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_s, x'_{T_s}, W_s) - T_s \cdot J(\theta^*, C^{\text{alg}}_s, T_s, x'_{T_s}, W_s)|$$

$$= \tilde{O}_T \left( \sum_{t=T_s}^{T_{s+1}-1} |C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x'_t) - C^{\text{alg}}_s(x'_t)| \right) + \tilde{O}_T(T_s \epsilon_s).$$

The proof of Lemma 16 can be found in Appendix E.13.

The control $C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x'_t)$ is safe for dynamics $\hat{\theta}_s$ and conditional on event $E$, $\|\hat{\theta}_s - \theta^*\|_\infty \le \tilde{O}_T(\nu_T) \le 1/\log(T)$ for sufficiently large $T$. The controller $C^{\text{alg}}_s$ is safe for dynamics $\theta^*$ for all $T$ steps conditional on event $E$ by definition of $E$. These together imply by Lemma 4 that, conditional on event $E$ and for sufficiently large $T$, for all $t \in [T_s, T_{s+1} - 1]$,

$$|x'_t|, |\hat{x}_t| \le 4 \log^2(T) \le B_x. \tag{57}$$

By Lemma 8 and 16, we have that conditional on event $E$, with probability $1 - o_T(1/T)$,

$$\sum_{s=0}^{s_e} T_s J(\theta^*, C^{\text{alg}}_s, T_s, x'_{T_s}, W_s) - \sum_{s=0}^{s_e} T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_s, \hat{x}_{T_s}, W_s)$$

$$= \tilde{O}_T(1) + \sum_{s=0}^{s_e} T_s J(\theta^*, C^{\text{alg}}_s, T_s, x'_{T_s}, W_s) - \sum_{s=0}^{s_e} T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_s, x'_{T_s}, W_s) \quad \text{Eq. (57), Lemma 8}$$

$$= \tilde{O}_T(1) + \sum_{s=0}^{s_e} \left( T_s J(\theta^*, C^{\text{alg}}_s, T_s, x'_{T_s}, W_s) - T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_s, x'_{T_s}, W_s) \right)$$

$$= \tilde{O}_T(1) + \tilde{O}_T \left( \sum_{s=0}^{s_e} \left( T_s \epsilon_s + \sum_{t=T_s}^{T_{s+1}-1} |C^{\text{alg}}_s(x'_t) - C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x'_t)| \right) \right) \quad \text{Lemma 16}$$

$$= \tilde{O}_T \left( \sum_{s=0}^{s_e} \epsilon_s T_s \right) + \tilde{O}_T \left( \sum_{s=0}^{s_e} \sum_{t=T_s}^{T_{s+1}-1} X^U_t \cdot \left| u^{\text{safeU}}_t - C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x'_t) \right| + X^L_t \cdot \left| u^{\text{safeL}}_t - C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x'_t) \right| \right).$$

We applied Lemma 8 for every $s \in [0 : s_e]$, so $\tilde{O}_T(1)$ times. Since Lemma 8 holds with probability $1 - o_T(1/T^{10})$, a union bound gives the first inequality holds with probability $1 - o_t(1/T^9)$. Another union bound combining this with the single application of Lemma 16 gives that the probability of the above result is $1 - o_T(1/T)$. The final line simplified using the fact that the two controls are equal if $X^L_t = X^U_t = 0$. $\square$

## E.8 Proof of Lemma 10 (Difference in Safety Controls)

*proof.* By symmetry, it is sufficient to show the first part of the lemma statement for $u^{\text{safeU}}_t$.

Because $C^{\text{alg}}$ is safe for dynamics $\theta^*$ under event $E$ and $E \subseteq E_1$, we have by Lemma 4 that under event $E$,

$$|x'_t| \le 4 \log^2(T). \tag{58}$$

Under event $E$ and for sufficiently large $T$, $\|\theta^* - \hat{\theta}_s\|_\infty \le \epsilon_s \le \frac{1}{\log(T)}$. This implies by construction of $u^{\text{safeU}}_t$ that under event $E$ and for sufficiently large $T$, $a^* x'_t + b^* u^{\text{safeU}}_t \le D_U$.

By Lemma 3, we also have that under event $E$ and for sufficiently large $T$, $u_t^{\text{safeU}} \geq u_t^{\text{safeL}}$. Therefore, by construction of $u_t^{\text{safeL}}$ we have that under event $E$ and for sufficiently large $T$, $a^* x_t' + b^* u_t^{\text{safeU}} \geq a^* x_t' + b^* u_t^{\text{safeL}} \geq D_{\text{L}}$. Together, this shows that $u_t^{\text{safeU}}$ is safe for dynamics $\theta^*$. By Lemma 4 and Equation (58), this gives that under event $E$ and for sufficiently large $T$,

$$|u_t^{\text{safeU}}| \leq B_x. \tag{59}$$

Because any control used by controller $C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}$ is safe for dynamics $\hat{\theta}_s$, by Lemma 4 we also have that under event $E$ for sufficiently large $T$,

$$\left| C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x_t') \right| \leq B_x. \tag{60}$$

Also, note that by Algorithm 2 Line 11, $u_t^{\text{safeU}}$ satisfies, for some $\theta$ such that $\|\theta - \hat{\theta}_s\|_\infty \leq \epsilon_s$,

$$ax_t' + bu_t^{\text{safeU}} = D_{\text{U}}. \tag{61}$$

Under event $E$, $\|\theta^* - \hat{\theta}_s\|_\infty \leq \epsilon_s$, which implies that $\|\theta^* - \theta\|_\infty \leq 2\epsilon_s \leq \tilde{O}_T(\nu_T) \leq 1/\log(T)$ for sufficiently large $T$. Therefore, applying Lemma 4 gives that under event $E$ and for sufficiently large $T$,

$$
\begin{aligned}
D_{\text{U}} &\geq a^* x_t' + b^* u_t^{\text{safeU}} & & u_t^{\text{safeU}} \text{ safe for } \theta^* \\
&\geq ax_t' + bu_t^{\text{safeU}} - |u_t^{\text{safeU}}| 2\epsilon_s - |x_t'| 2\epsilon_s & & \|\theta^* - \theta\|_\infty \leq 2\epsilon_s \\
&\geq D_{\text{U}} - 4B_x \epsilon_s. & & \text{Equations (58),(59), and (61)} 
\end{aligned}
\tag{62}
$$

If $u_t^{\text{safeU}} \leq C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x_t')$, then there must exist some $\theta$ such that $\|\hat{\theta}_s - \theta\|_\infty \leq \epsilon_s$ and

$$ax_t' + bC^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x_t') \geq D_{\text{U}}. \tag{63}$$

Under event $E$, $\|\theta^* - \theta\|_\infty \leq 2\epsilon_s \leq \tilde{O}_T(\nu_T) \leq 1/\log(T)$ for sufficiently large $T$, therefore under event $E$ and for sufficiently large $T$,

$$
\begin{aligned}
&a^* x_t' + b^* C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x_t') \\
&\geq ax_t' + bC^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x_t') - 2\epsilon_s |x_t'| - 2\epsilon_s \left| C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x_t') \right| \\
&\geq D_{\text{U}} - 4B_x \epsilon_s. & \text{Equations (58),(60), and (63)}
\end{aligned}
\tag{64}
$$

Finally, because $C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x_t')$ is safe for dynamics $\hat{\theta}_s$,

$$\hat{a}_s x_t' + \hat{b}_s C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x_t') \leq D_{\text{U}}. \tag{65}$$

Using that under event $E$, $\|\theta^* - \hat{\theta}_s\|_\infty \leq \epsilon_s \leq \tilde{O}_T(\nu_T) \leq 1/\log(T)$ for sufficiently large $T$, Equations 58, 60, and (65) imply that under event $E$ and for sufficiently large $T$,

$$a^* x_t' + b^* C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x_t') \leq D_{\text{U}} + 2B_x \epsilon_s. \tag{66}$$

47

Combining Equations (64) and (66), if $u_t^{\text{safeU}} \leq C_{K_{\text{opt}}(\hat{\theta}_s,T_s)}^{\hat{\theta}_s}(x_t')$ then under event $E$ and for sufficiently large $T$,

$$D_{\text{U}} - 4B_x\epsilon_s \leq a^*x_t' + b^*C_{K_{\text{opt}}(\hat{\theta}_s,T_s)}^{\hat{\theta}_s}(x_t') \leq D_{\text{U}} + 2B_x\epsilon_s.$$

Combining this with Equation (62) gives that under event $E$ and for sufficiently large $T$,

$$|(a^*x_t' + b^*u_t^{\text{safeU}}) - (a^*x_t' + b^*C_{K_{\text{opt}}(\hat{\theta}_s,T_s)}^{\hat{\theta}_s}(x_t'))| = 6B_x\epsilon_s.$$

This implies the desired result that under event $E$ and for sufficiently large $T$,

$$|u_t^{\text{safeU}} - C_{K_{\text{opt}}(\hat{\theta}_s,T_s)}^{\hat{\theta}_s}(x_t')| = 6B_x\epsilon_s/b^*.$$

$\square$

## E.9   Proof of Lemma 15 (McDiarmid's Condition)

*proof.* First, we will construct the event $E_s^{\text{M}}$. Define

$$E_s^{\text{M}} = \{\forall t \in [T_s : T_{s+1} - 1], |w_t| \leq \log^2(T)\} \cap \bigcap_{i=T_s}^{T_{s+1}-1} E_{\text{A8}}\left(C_{K_{\text{opt}}(\hat{\theta}_s,T_s)}^{\hat{\theta}_s}, \{w_t\}_{t=i}^{T_{s+1}-1}\right).$$

Note because $\mathbb{P}(\{\forall t \in [T_s : T_{s+1} - 1], |w_t| \leq \log^2(T)\}) \geq \mathbb{P}(E_1) = 1 - o_T(1/T^{10})$ and because under event $E_2^s$, $\mathbb{P}\left(E_{\text{A8}}\left(C_{K_{\text{opt}}(\hat{\theta}_s,T_s)}^{\hat{\theta}_s}, \{w_t\}_{t=i}^{T_{s+1}-1}\right) \mid \hat{\theta}_s\right) = 1 - o_T(1/T^{10})$ we have by a union bound that $\mathbb{P}(E_s^{\text{M}} \mid \hat{\theta}_s) = 1 - o_T(1/T^9)$. Suppose $E_s^{\text{M}}$ holds for $W_s$ and $W_s'$. For $i \in [T_s, T_{s+1}]$, define $W^i$ as follows.

$$W^i = \{w_{T_s}, w_{T_s+1}, ..., w_{i-1}, w_i', w_{i+1}', w_{i+2}', ...w_{T_{s+1}-1}'\}.$$

In other words, $W^i$ includes noise $w_t$ for $t < i$ and includes $w_t'$ for $t \geq i$. For $i \in [T_s, T_{s+1} - 1]$, we will first bound

$$\left|f_{\hat{\theta}_s}(W^i) - f_{\hat{\theta}_s}(W^{i+1})\right|.$$

First, note that if $w_i = w_i'$, then $W^i = W^{i+1}$ and therefore $f_{\hat{\theta}_s}(W^i) = f_{\hat{\theta}_s}(W^{i+1})$. Now, assume $w_i \neq w_i'$. Let $x_0^i, .., x_{T_s}^i$ be the series of positions when the noise random variables are $W^i$, $x_0^i = 0$, and the controller used is $C_{K_{\text{opt}}(\hat{\theta}_s,T_s)}^{\hat{\theta}_s}$. Conditional on $E_2^s$, $\|\hat{\theta}_s - \theta^*\|_\infty \leq \tilde{O}(\nu_T) \leq 1/\log(T)$ for sufficiently large $T$. Because $E_s^{\text{M}}$ holds for $W_s, W_s'$, we have that $E_1$ holds for $W^i$ for all $i$. Therefore by Lemma 4 for sufficiently large $T$, $|x_t^i| \leq 4\log^2(T)$ for all $i, t$. For any $t \leq i$, $x_t^i = x_t^{i+1}$. Therefore, the difference in the two trajectories $\{x_t^i\}$ and $\{x_t^{i+1}\}$ only occurs at and after time $i + 1$. The first difference occurs at time $i + 1$ when $x_{i+1}^i = x_{i+1}^{i+1} - w_i + w_i'$. For the next $T_{s+1} - i - 1$ steps, the difference in cost of the two

48

trajectories $\{x_t^i\}$ and $\{x_t^{i+1}\}$ is

$$(T_{s+1} - i - 1)J(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_{s+1} - i - 1, x_{i+1}^{i+1}, \{w'_t\}_{t=i+1}^{T_{s+1}-1})$$

$$- (T_{s+1} - i - 1)J(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_{s+1} - i - 1, x_{i+1}^i, \{w'_t\}_{t=i+1}^{T_{s+1}-1})$$

$$= \tilde{O}_T\left(|x_{i+1}^{i+1} - x_{i+1}^i| + |\hat{\theta}_s - \theta^*|\right) \qquad\qquad \text{Lemma 8, } |x_t^i| \le 4\log^2(T)$$

$$= \tilde{O}_T\left(|x_{i+1}^{i+1} - x_{i+1}^{i+1} + w_i - w'_i| + \nu_T\right) \qquad\qquad \text{Event } E_2^s, \ x_{i+1}^i = x_{i+1}^{i+1} - w_i + w'_i$$

$$= \tilde{O}_T\left(|w_i - w'_i| + \nu_T\right)$$

$$= \tilde{O}_T\left(2\log^2(T) + \nu_T\right) \qquad\qquad W, W' \text{ satisfy event } E_s^{\text{M}}$$

$$= \tilde{O}_T(1). \tag{67}$$

We have therefore shown that for some $c = \tilde{O}_T(1)$,

$$|f_{\hat{\theta}_s}(W^i) - f_{\hat{\theta}_s}(W^{i+1})| \le c.$$

Because $W_s = W^{T_{s+1}}$ and $W'_s = W^{T_s}$, we have by the triangle inequality that

$$|f_{\hat{\theta}_s}(W_s) - f_{\hat{\theta}_s}(W'_s)| = |f_{\hat{\theta}_s}(W^{T_{s+1}}) - f_{\hat{\theta}_s}(W^{T_s})|$$

$$\le \sum_{i=T_s}^{T_{s+1}-1} |f_{\hat{\theta}_s}(W^i) - f_{\hat{\theta}_s}(W^{i+1})|$$

$$= \sum_{i=T_s, w_i \ne w'_i}^{T_{s+1}-1} |f_{\hat{\theta}_s}(W^i) - f_{\hat{\theta}_s}(W^{i+1})|$$

$$\le \sum_{i=T_s, w_i \ne w'_i}^{T_{s+1}-1} c.$$

$\square$

## E.10 Proof of Lemma 11

*proof.* Define $E^* = \{|x|, |y| \le 4\log^2(T)\} \cap E_{\text{A8}}(C_K^\theta, W')$. By assumption of the lemma, we have that $\mathbb{P}(|x| \le 4\log^2(T)) = 1 - o_T(1/T^{11})$ and $\mathbb{P}(|y| \le 4\log^2(T)) = 1 - o_T(1/T^{11})$. Because $\|\theta - \theta^*\|_\infty \le \epsilon_{\text{A8}}$, $\mathbb{P}(E_{\text{A8}}(C_K^\theta, W')) = 1 - o_T(1/T^{10})$. Therefore, by a union bound we have that $\mathbb{P}(E^*) = 1 - o_T(1/T^{10})$. By the Law of Total Expectation,

$$\mathbb{E}[|t \cdot J(\theta^*, C_K^\theta, t, x, W') - t \cdot J(\theta^*, C_K^\theta, t, y, W')|]$$

$$= \mathbb{E}\left[|t \cdot J(\theta^*, C_K^\theta, t, x, W') - t \cdot J(\theta^*, C_K^\theta, t, y, W')| \,\big|\, E^*\right] \mathbb{P}(E^*)$$

$$\quad + \mathbb{E}\left[|t \cdot J(\theta^*, C_K^\theta, t, x, W') - t \cdot J(\theta^*, C_K^\theta, t, y, W')| \,\big|\, \neg E^*\right] \mathbb{P}(\neg E^*)$$

$$= \mathbb{E}\left[\tilde{O}_T\left(|x - y| + \epsilon\right) \,\big|\, E^*\right] \mathbb{P}(E^*)$$

$$\quad + \mathbb{E}\left[|t \cdot J(\theta^*, C_K^\theta, t, x, W') - t \cdot J(\theta^*, C_K^\theta, t, y, W')| \,\big|\, \neg E^*\right] \mathbb{P}(\neg E^*) \qquad \text{Lemma 8}$$

$$= \tilde{O}_T\left(\mathbb{E}[|x - y| \mid E^*]\,\mathbb{P}(E^*) + \epsilon\right)$$

$$\quad + \mathbb{E}\left[|t \cdot J(\theta^*, C_K^\theta, t, x, W') - t \cdot J(\theta^*, C_K^\theta, t, y, W')| \,\big|\, \neg E^*\right] \mathbb{P}(\neg E^*)$$

$$= \tilde{O}_T\left(\mathbb{E}[|x - y|] + \epsilon\right) + \mathbb{E}\left[|t \cdot J(\theta^*, C_K^\theta, t, x, W') - t \cdot J(\theta^*, C_K^\theta, t, y, W')| \,\big|\, \neg E^*\right] \mathbb{P}(\neg E^*) \quad \text{LoTE}$$

Therefore, all we must show is that

$$\mathbb{E}\left[\left|t \cdot J(\theta^*, C_K^\theta, t, x, W') - t \cdot J(\theta^*, C_K^\theta, t, y, W')\right| \mid \neg E^*\right] \mathbb{P}(\neg E^*) = \tilde{O}_T(T^{-2}).$$

Define $w_m = \max_{w \in W'} |w|$. By Lemma 12, we can bound the position and controls at every time step in terms of $w_m$ to get that

$$\begin{aligned}
&\left|t \cdot J(\theta^*, C_K^\theta, t, x, W') - t \cdot J(\theta^*, C_K^\theta, t, y, W')\right| \\
&= T(q+r)O_T\left((w_m + x + \|D\|_\infty)^2 + (w_m + y + \|D\|_\infty)^2\right) \qquad \text{Triangle Inequality, Lemma 12} \\
&= O_T\left(T\left((w_m + x + \|D\|_\infty)^2 + (w_m + y + \|D\|_\infty)^2\right)\right) \\
&= \tilde{O}_T\left(T\left(w_m^2 + w_m|x| + |x|^2 + w_m|y| + |y|^2 + w_m + |x| + |y| + 1\right)\right). \quad \text{Assum 3 } (\|D\|_\infty \le \log^2(T))
\end{aligned}$$

Therefore, we have that

$$\begin{aligned}
&\mathbb{E}\left[\left|t \cdot J(\theta^*, C_K^\theta, t, x, W') - t \cdot J(\theta^*, C_K^\theta, t, y, W')\right| \mid \neg E^*\right] \mathbb{P}(\neg E^*) \\
&= \tilde{O}_T\Big(T\big(\mathbb{E}[w_m^2 \mid \neg E^*]\mathbb{P}(\neg E^*) + \mathbb{E}[w_m \mid \neg E^*]\mathbb{P}(\neg E^*) + \mathbb{E}[|y|w_m \mid \neg E^*]\mathbb{P}(\neg E^*) + \mathbb{E}[|y|^2 \mid \neg E^*]\mathbb{P}(\neg E^*) \\
&\quad + \mathbb{E}[|x|w_m \mid \neg E^*]\mathbb{P}(\neg E^*) + \mathbb{E}[|x|^2 \mid \neg E^*]\mathbb{P}(\neg E^*) + \mathbb{E}[|x| \mid \neg E^*]\mathbb{P}(\neg E^*) + \mathbb{E}[|y| \mid \neg E^*]\mathbb{P}(\neg E^*) + \mathbb{P}(\neg E^*)\big)\Big).
\end{aligned}$$

$$(68)$$

Therefore, it is sufficient to show that $\mathbb{E}[w_m \mid \neg E^*]\mathbb{P}(\neg E^*)$, $\mathbb{E}[w_m^2 \mid \neg E^*]\mathbb{P}(\neg E^*)$, $\mathbb{E}[|x| \mid \neg E^*]\mathbb{P}(\neg E^*)$, $\mathbb{E}[x^2 \mid \neg E^*]\mathbb{P}(\neg E^*)$, $\mathbb{E}[|y| \mid \neg E^*]\mathbb{P}(\neg E^*)$, $\mathbb{E}[y^2 \mid \neg E^*]\mathbb{P}(\neg E^*)$, $\mathbb{E}[|x|w_m \mid \neg E^*]\mathbb{P}(\neg E^*)$, $\mathbb{E}[|y|w_m \mid \neg E^*]\mathbb{P}(\neg E^*)$ are all $\tilde{O}_T(\frac{1}{T^3})$. We will use the following probability result.

**Lemma 17.** *Suppose $X$ is a non-negative random variable. Then for any $L \ge 0$ and any event $E$, we have that*

$$\mathbb{E}\left[X \mid E\right]\mathbb{P}(E) \le \mathbb{P}(E)L + \mathbb{P}(X \ge L)\mathbb{E}\left[X \mid X \ge L\right]$$

*proof.* For any events $A, B$ such that $A \subseteq B$, we have that

$$\begin{aligned}
\mathbb{E}[X \mid B]\mathbb{P}(B) &= \mathbb{E}[X \mid A, B]\mathbb{P}(A \mid B)\mathbb{P}(B) + \mathbb{E}[X \mid \neg A, B]\mathbb{P}(\neg A \mid B)\mathbb{P}(B) \\
&= \mathbb{E}[X \mid A]\mathbb{P}(A) + \mathbb{E}[X \mid \neg A, B]\mathbb{P}(\neg A \mid B)\mathbb{P}(B) \qquad A \subseteq B \\
&\ge \mathbb{E}[X \mid A]\mathbb{P}(A).
\end{aligned}$$

$$(69)$$

Therefore, we can conclude that

$$\begin{aligned}
&\mathbb{E}\left[X \mid E\right]\mathbb{P}(E) \\
&= \mathbb{E}\left[X \mid E, X \le L\right]\mathbb{P}(X \le L \mid E)\mathbb{P}(E) + \mathbb{E}\left[X \mid E, X \ge L\right]\mathbb{P}(X \ge L \mid E)\mathbb{P}(E) \\
&\le \mathbb{P}(E)L + \mathbb{E}\left[X \mid E, X \ge L\right]\mathbb{P}(X \ge L \mid E)\mathbb{P}(E) \\
&\le \mathbb{P}(E)L + \mathbb{E}\left[X \mid E, X \ge L\right]\mathbb{P}(E, X \ge L) \\
&\le \mathbb{P}(E)L + \mathbb{E}\left[X \mid X \ge L\right]\mathbb{P}(X \ge L). \qquad \text{Eq (69)}
\end{aligned}$$

$\square$

Now, note that by the assumption on $x$ and definition of $E^*$ (where $L$ is from the lemma statement),

$$\mathbb{E}[x^2 \mid \neg E^*]\mathbb{P}(\neg E^*) \leq \mathbb{P}(\neg E^*)L^2 + \mathbb{P}(|x| \geq L)\,\mathbb{E}[|x|^2 \mid |x| \geq L] \qquad \text{Lemma 17}$$

$$= \tilde{O}_T\left(\frac{1}{T^{10}}\right) + \tilde{O}_T\left(\frac{1}{T^{10}}\right)$$

$$= \tilde{O}_T\left(\frac{1}{T^{10}}\right).$$

This also implies by the Cauchy–Schwarz inequality that

$$\mathbb{E}[|x| \mid \neg E^*]\mathbb{P}(\neg E^*) \leq \sqrt{\mathbb{E}[x^2 \mid \neg E^*]}\mathbb{P}(\neg E^*)$$

$$= \sqrt{\mathbb{E}[x^2 \mid \neg E^*]\mathbb{P}(\neg E^*)}\sqrt{\mathbb{P}(\neg E^*)}$$

$$= \tilde{O}_T\left(\frac{1}{T^5}\right).$$

By Lemma 13, because $\mathbb{P}(E^*) = 1 - o_T(1/T^{11})$ we have that

$$\mathbb{E}[w_m^2 \mid \neg E^*]\mathbb{P}(\neg E^*) = \tilde{O}_T\left(\frac{1}{T^{10}}\right).$$

Once again, by the Cauchy-Schwarz inequality this implies that $\mathbb{E}[w_m \mid \neg E^*] = \tilde{O}_T\left(\frac{1}{T^5}\right)$.

By the subgaussian assumption on $\mathcal{D}$ and a union bound, we have that

$$\mathbb{P}(w_m \geq \log^3(T)) \leq \sum_{w \in W'} \mathbb{P}(|w| \geq \log^3(T))$$

$$\leq t \cdot 2e^{-\Omega_T(\log^6(T))}$$

$$\leq o_T(1/T^{11}). \tag{70}$$

Finally, we have by the independence of $x$ and $w_m$ and the assumption on $x$ that

$\mathbb{E}[|x|w_m \mid \neg E^*]\mathbb{P}(\neg E^*)$

$\leq \mathbb{P}(\neg E^*)L\log^3(T)$

$\qquad + \mathbb{P}(|x| \geq L, w_m \geq \log^3(T))\,\mathbb{E}\left[|x|w_m \mid |x| \geq L, w_m \geq \log^3(T)\right]$

$\qquad + \mathbb{P}(|x| \leq L, w_m \geq \log^3(T))\,\mathbb{E}[|x|w_m \mid |x| \leq L, w_m \geq \log^3(T)]$

$\qquad + \mathbb{P}(|x| \geq L, w_m \leq \log^3(T))\,\mathbb{E}[|x|w_m \mid |x| \geq L, w_m \leq \log^3(T)] \qquad \text{Lemma 17}$

$\leq \mathbb{P}(\neg E^*)L\log^3(T)$

$\qquad + \mathbb{P}(|x| \geq L)\mathbb{P}(w_m \geq \log^3(T))\,\mathbb{E}\left[|x| \mid |x| \geq L\right]\mathbb{E}\left[w_m \mid w_m \geq \log(T)\right]$

$\qquad + L\mathbb{P}(w_m \geq \log^3(T))\,\mathbb{E}[w_m \mid w_m \geq \log^3(T)]$

$\qquad + \log^3(T)\mathbb{P}(|x| \geq L)\,\mathbb{E}\left[|x| \mid |x| \geq L\right] \qquad\qquad \text{[Ind of } x \text{ and } w_m]$

$= \tilde{O}_T\left(\frac{1}{T^{10}}\right) + \tilde{O}_T\left(\frac{1}{T^{20}}\right) + \tilde{O}_T\left(\frac{1}{T^{10}}\right) + \tilde{O}_T\left(\frac{1}{T^{10}}\right) \qquad \text{[Def of } E^*, \text{ Lemma 13, Eq (70), Assum on } x]$

$= \tilde{O}_T\left(\frac{1}{T^{10}}\right).$

Note that by symmetry, all of the above results also hold for $y$. Returning to Equation (68), we have that

$$\mathbb{E}\left[\left|t \cdot J(\theta^*, C_K^\theta, t, x, W') - t \cdot J(\theta^*, C_K^\theta, t, y, W')\right| \mid \neg E^*\right] \mathbb{P}(\neg E^*) = \tilde{O}_T\left(\frac{1}{T^2}\right).$$

This completes the proof that

$$\mathbb{E}\left[\left|t \cdot J(\theta^*, C_K^\theta, t, x, W') - t \cdot J(\theta^*, C_K^\theta, t, y, W')\right|\right] = \tilde{O}_T\left(\mathbb{E}[|x - y|] + \epsilon + \frac{1}{T^2}\right).$$

$\square$

## E.11   Proof of Lemma 12

*proof.* Define $\gamma_T = \max_{t \leq T-1} \|\theta_t - \theta^*\|_\infty \leq \frac{1}{\log(T)}$. Because the control at time $t$ is safe for dynamics $\theta_t$, we have $D_\mathrm{L} \leq a_t x_t + b_t u_t \leq D_\mathrm{U}$ for all $t$. By the triangle inequality,

$$|x_{t+1}| = |a^* x_t + b^* u_t + w_t| \leq |w_t| + \|D\|_\infty + \gamma_T(|x_t| + |u_t|).$$

As in Equation (52),

$$|u_t| \leq \frac{\|D\|_\infty + a^* |x_t| + \gamma_T |x_t|}{b^* - \gamma_T} = \frac{\|D\|_\infty + (a^* + \gamma_T)|x_t|}{b^* - \gamma_T}.$$

For sufficiently large $T$, $\gamma_T \leq b^*/2$, and therefore for sufficiently large $T$,

$$|x_{t+1}| \leq |w_t| + \|D\|_\infty + \gamma_T\left(|x_t| + \frac{\|D\|_\infty + a^*|x_t| + \gamma_T|x_t|}{b^* - \gamma_T}\right) = O_T(|w_t| + \|D\|_\infty + \gamma_T|x_t|).$$

Using $x_0 = x$ as the base-case, this recursive relationship implies that for all $t$,

$$\begin{aligned}
|x_t| &\leq O_T\left(\sum_{i=0}^{t-1}(|w_i| + \|D\|_\infty + |x|)\gamma_T^{t-1-i}\right) \\
&\leq O_T\left(\left(\max_{i \leq t-1}|w_i| + \|D\|_\infty + |x|\right)\sum_{i=0}^{t-1}\gamma_T^i\right) \\
&\leq O_T\left(\left(\max_{i \leq t-1}|w_i| + \|D\|_\infty + |x|\right)\sum_{i=0}^{t-1}\left(\frac{1}{\log(T)}\right)^i\right) \\
&= O_T\left(\left(\max_{i \leq t-1}|w_i| + \|D\|_\infty + |x|\right)\frac{1}{1 - \frac{1}{\log(T)}}\right).
\end{aligned}$$

This implies that for sufficiently large $T$, $|x_t| = O_T(\max_{i \leq t-1}|w_i| + \|D\|_\infty + |x|)$ and

$$|u_t| \leq \frac{\|D\|_\infty + (a^* + \gamma_T)O_T(\max_{i \leq t-1}|w_i| + \|D\|_\infty + |x|)}{b^* - \gamma_T} = O_T(\max_{i \leq t-1}|w_i| + \|D\|_\infty + |x|),$$

which are exactly the desired bounds.   $\square$

## E.12 Proof of Lemma 13

*proof.* Define $w_m = \max_{i \leq t} |w_i|$. Because $w_t$ is sub-Gaussian, there exists $\alpha > 0$ such that for any $w \geq 0$, $\mathbb{P}(|w_t| \geq w) \leq 2e^{-w^2/(2\alpha)}$. Therefore, we have for any $w \geq 0$,

$$
\begin{aligned}
\mathbb{P}(w_m^2 \geq w) &= 1 - \mathbb{P}\left(\forall i \leq t, |w_i| \leq \sqrt{w}\right) \\
&\leq 1 - \left(1 - 2e^{-w/(2\alpha)}\right)^t \\
&= O_T(te^{-w/(2\alpha)}).
\end{aligned}
$$

This implies by the Law of Total Expectation that

$$
\begin{aligned}
\mathbb{E}[w_m^2 \mid \neg F]\mathbb{P}(\neg F) &\leq \mathbb{P}(\neg F)\log^6(T) + \mathbb{P}(w_m \geq \log^3(T))\,\mathbb{E}[w_m^2 \mid w_m \geq \log^3(T)] \quad \text{Lemma 17} \\
&= o_T\left(\frac{1}{T^{10}}\right) + \int_{\log^6(T)}^{\infty} \mathbb{P}(w_m^2 \geq w)dw \\
&= o_T\left(\frac{1}{T^{10}}\right) + O_T\left(\int_{\log^6(T)}^{\infty} te^{-w/(2\alpha)}dw\right) \\
&= o_T\left(\frac{1}{T^{10}}\right) + O_T\left(2t\alpha e^{-\log^6(T)/(2\alpha)}\right) \\
&= o_T\left(\frac{1}{T^{10}}\right).
\end{aligned}
$$

$\square$

## E.13 Proof of Lemma 16

*proof.* Fix a value of $s$. For $i \in [0 : T_s]$, define the controller $C_t^i$ as the controller that at time $t < i$ uses controller $C_s^{\text{alg}}$ and at time $t \geq i$ uses controller $C_{K_{\text{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}$. We will compare the cost of controller $C_t^i$ versus controller $C_t^{i+1}$ over $T_s$ steps starting at position $x'_{T_s}$. Note that the cost of the first $i$ steps is the same, as $C_t^i = C_t^{i+1} = C_s^{\text{alg}}$ for $t < i$. Therefore

$$
|i \cdot J(\theta^*, C_t^{i+1}, i, x'_{T_s}, \{w_j\}_{j=T_s}^{T_s+i-1}) - i \cdot J(\theta^*, C_t^i, i, x'_{T_s}, \{w_j\}_{j=T_s}^{T_s+i-1})| = 0.
$$

The position at time $i$ when using either $C_t^i$ or $C_t^{i+1}$ is $x'_{T_s+i}$. Conditional on event $E$ and for sufficiently large $T$, by Lemma 4 we have that $|C_s^{\text{alg}}(x'_{T_s+i})|, |C_{K_{\text{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}(x'_{T_s+i})| \leq B_x$. Therefore conditional on event $E$ and for sufficiently large $T$,

$$
\begin{aligned}
&r\left(C_s^{\text{alg}}(x'_{T_s+i})^2 - C_{K_{\text{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}(x'_{T_s+i})^2\right) \\
&\leq 2r|C_s^{\text{alg}}(x'_{T_s+i}) - C_{K_{\text{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}(x'_{T_s+i})| + r\left(C_s^{\text{alg}}(x'_{T_s+i}) - C_{K_{\text{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}(x'_{T_s+i})\right)^2 \\
&\leq r(2 + 2B_x)|C_s^{\text{alg}}(x'_{T_s+i}) - C_{K_{\text{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}(x'_{T_s+i})|. \quad (71)
\end{aligned}
$$

The difference in the next position when at position $x'_{T_s+i}$ and using control $C_s^{\text{alg}}(x'_{T_s+i})$ versus $C^{\hat\theta_s}_{K_{\text{opt}}(\hat\theta_s,T_s)}(x'_{T_s+i})$ is

$$\left| a^* x'_{T_s+i} + b^* C_s^{\text{alg}}(x'_{T_s+i}) + w_{T_s+i} - \left( a^* x'_{T_s+i} + b^* C^{\hat\theta_s}_{K_{\text{opt}}(\hat\theta_s,T_s)}(x'_{T_s+i}) + w_{T_s+i} \right) \right|$$
$$= b^* \left| C_s^{\text{alg}}(x'_{T_s+i}) - C^{\hat\theta_s}_{K_{\text{opt}}(\hat\theta_s,T_s)}(x'_{T_s+i}) \right|. \tag{72}$$

Under event $E$, the controls $C^{\hat\theta_s}_{K_{\text{opt}}(\hat\theta_s,T_s)}(x'_{T_s+i})$ and $C_s^{\text{alg}}(x'_{T_s+i})$ are safe for dynamics $\hat\theta_s$ and $\theta^*$, respectively and $\|\theta^* - \hat\theta_s\|_\infty \le \epsilon_s \le \tilde O_T(\nu_T)$. Therefore, by Lemma 4, conditional on event $E$ and for sufficiently large $T$, we have that $|x'_{T_s+i}| \le 4\log^2(T)$ and that

$$|a^* x'_{T_s+i} + b^* C_s^{\text{alg}}(x'_{T_s+i}) + w_{T_s+i}|, |a^* x'_{T_s+i} + b^* C^{\hat\theta_s}_{K_{\text{opt}}(\hat\theta_s,T_s)}(x'_{T_s+i}) + w_{T_s+i}| \le 4\log^2(T).$$

Conditional on event $E$ and for sufficiently large $T$, we therefore have by Lemma 8 that conditional on $E_{\text{A8}}(C^{\hat\theta_s}_{K_{\text{opt}}(\hat\theta_s,T_s)}, \{w_j\}_{j=T_s+i+1}^{T_{s+1}-1})$, we can bound the difference in future cost of the next $T_s - i - 1$ steps starting at time $T_s + i + 1$ when using control $C_s^{\text{alg}}(x'_{T_s+i})$ versus $C^{\hat\theta_s}_{K_{\text{opt}}(\hat\theta_s,T_s)}(x'_{T_s+i})$ as follows.

$$(T_s - i - 1)\left| J(\theta^*, C^{\hat\theta_s}_{K_{\text{opt}}(\hat\theta_s,T_s)}, T_s - i - 1, a^* x'_{T_s+i} + b^* C_s^{\text{alg}}(x'_{T_s+i}) + w_{T_s+i}, \{w_j\}_{j=T_s+i+1}^{T_{s+1}-1}) \right.$$
$$\left. - J(\theta^*, C^{\hat\theta_s}_{K_{\text{opt}}(\hat\theta_s,T_s)}, T_s - i - 1, a^* x'_{T_s+i} + b^* C^{\hat\theta_s}_{K_{\text{opt}}(\hat\theta_s,T_s)}(x'_{T_s+i}) + w_{T_s+i}, \{w_j\}_{j=T_s+i+1}^{T_{s+1}-1}) \right|$$
$$= \tilde O_T\left( b^* |C_s^{\text{alg}}(x'_{T_s+i}) - C^{\hat\theta_s}_{K_{\text{opt}}(\hat\theta_s,T_s)}(x'_{T_s+i})| + \epsilon_s \right). \qquad\text{[Eq (72), Lemma 8]} \tag{73}$$

Therefore, the difference in total cost between $C_t^i$ and $C_t^{i+1}$ conditional on event $E$ with probability $\mathbb{P}(E_{\text{A8}}(C^{\hat\theta_s}_{K_{\text{opt}}(\hat\theta_s T_s)}, \{w_j\}_{j=T_s+i+1}^{T_{s+1}-1})) = 1 - o_T(1/T^{10})$ is

$$|T_s \cdot J(\theta^*, C_t^{i+1}, T_s, x'_{T_s}, W_s) - T_s \cdot J(\theta^*, C_t^i, T_s, x'_{T_s}, W_s)$$
$$\le \left| i \cdot J(\theta^*, C_t^{i+1}, i, x'_{T_s}, \{w_j\}_{j=T_s}^{T_s+i}) - i \cdot J(\theta^*, C_t^i, i, x'_{T_s}, \{w_j\}_{j=T_s}^{T_s+i}) \right| + r\left( C_s^{\text{alg}}(x'_{T_s+i})^2 - C^{\hat\theta_s}_{K_{\text{opt}}(\hat\theta_s,T_s)}(x'_{T_s+i})^2 \right)$$
$$+ (T_s - i - 1)\left| J(\theta^*, C^{\hat\theta_s}_{K_{\text{opt}}(\hat\theta_s,T_s)}, T_s - i - 1, a^* x'_{T_s+i} + b^* C_s^{\text{alg}}(x'_{T_s+i}) + w_{T_s+i}, \{w_j\}_{j=T_s+i+1}^{T_{s+1}-1}) \right.$$
$$\left. - J(\theta^*, C^{\hat\theta_s}_{K_{\text{opt}}(\hat\theta_s,T_s)}, T_s - i - 1, a^* x'_{T_s+i} + b^* C^{\hat\theta_s}_{K_{\text{opt}}(\hat\theta_s,T_s)}(x'_{T_s+i}) + w_{T_s+i}, \{w_j\}_{j=T_s+i+1}^{T_{s+1}-1}) \right|$$
$$\le 0 + r(2 + 2B_x)|C_s^{\text{alg}}(x'_{T_s+i}) - C^{\hat\theta_s}_{K_{\text{opt}}(\hat\theta_s,T_s)}(x'_{T_s+i})|$$
$$+ \tilde O_T\left( b^* |C_s^{\text{alg}}(x'_{T_s+i}) - C^{\hat\theta_s}_{K_{\text{opt}}(\hat\theta_s,T_s)}(x'_{T_s+i})| + \epsilon_s \right) \quad \text{Eq. (71), (73)}$$
$$= \tilde O_T\left( |C_s^{\text{alg}}(x'_{T_s+i}) - C^{\hat\theta_s}_{K_{\text{opt}}(\hat\theta_s,T_s)}(x'_{T_s+i})| + \epsilon_s \right). \tag{74}$$

We can use Equation (74) for all $i \in [0 : T_s - 1]$, the triangle inequality, and a union bound to get that conditional on event $E$, with probability $1 - o_T(1/T^9)$

$$|T_s \cdot J(\theta^*, C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}, T_s, x'_{T_s}, W_s) - T_s \cdot J(\theta^*, C^{\text{alg}}_s, T_s, x'_{T_s}, W_s)|$$

$$\leq \sum_{i=0}^{T_s-1} |T_s \cdot J(\theta^*, C^{i+1}_t, T_s, x'_{T_s}, W_s) - T_s \cdot J(\theta^*, C^i_t, T_s, x'_{T_s}, W_s)|$$

$$= \tilde{O}_T \left( \sum_{i=0}^{T_s-1} |C^{\text{alg}}_s(x'_{T_s+i}) - C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x'_{T_s+i})| + T_s \epsilon_s \right). \tag{75}$$

The above was for a fixed value of $s$. Taking a union bound over all $s \in [0 : s_e]$, we have that with conditional probability $1 - o_T(1/T)$ given event $E$, the desired result holds for all $s \in [0 : s_e]$. $\qquad \square$

# F  Proofs of Sufficiently Large Noise Case

## F.1  Proof of Theorem 1

First, we present the algorithm which is used to prove Theorem 1.

---

**Algorithm 3** Safe LQR for Large Noise

---

*proof.* **Input:** $D, \mathcal{D}, \Theta, C^{\text{init}}, \{\mathcal{C}^\theta\}_{\theta \in \Theta}, T, \lambda$

1:  $\nu_T \leftarrow T^{-1/4}$
2:  **for** $t \leftarrow 0$ to $\frac{1}{\nu_T^2} - 1$ **do**          ▷ Safe warm-up exploration phase
3:      $\phi_t \sim \text{Rademacher}(0.5)$
4:      Use control $u_t = C^{\text{init}}(x_t) + \frac{\phi_t}{\log(T)}$
5:  **for** $s \leftarrow 0$ to $\log_2(T\nu_T^2) - 1$ **do**          ▷ Safe certainty equivalence phase
6:      $T_s \leftarrow \frac{2^s}{\nu_T^2}$
7:      $\epsilon_s \leftarrow B_{T_s} \sqrt{\frac{\max(V_{T_s}^{22}, V_{T_s}^{11})}{V_{T_s}^{11} V_{T_s}^{22} - (V_{T_s}^{12})^2}}$
8:      $\hat{\theta}_s^{\text{pre}} \leftarrow (Z_{T_s}^\top Z_{T_s} + \lambda I)^{-1} Z_{T_s}^\top X_{T_s}$
9:      $\hat{\theta}_s \leftarrow \arg\min_{\|\theta' - \hat{\theta}_s^{\text{pre}}\|_\infty \leq \epsilon_s} \min_{\|\theta - \hat{\theta}_s^{\text{pre}}\|_\infty \leq \epsilon_s} P(\theta, K_{\text{opt}}(\theta', T_s), D_{\text{U}})$
10:     $C^{\text{alg}}_s \leftarrow C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}$
11:     **for** $t \leftarrow T_s$ to $2T_s - 1$ **do**
12:         $u_t^{\text{safeU}} \leftarrow \max \left\{ u : \max_{\|\theta - \hat{\theta}_s^{\text{pre}}\|_\infty \leq \epsilon_s} ax_t + bu \leq D_{\text{U}} \right\}$
13:         $u_t^{\text{safeL}} \leftarrow \min \left\{ u : \min_{\|\theta - \hat{\theta}_s^{\text{pre}}\|_\infty \leq \epsilon_s} ax_t + bu \geq D_{\text{L}} \right\}$
14:         Use control $u_t = \max \left( \min \left( C^{\text{alg}}_s(x_t), u_t^{\text{safeU}} \right), u_t^{\text{safeL}} \right)$

---

Importantly, we note that Algorithm 3 fundamentally only differs from Algorithm 2 in two ways. The first is that $\nu_T$ changes from $T^{-1/3}$ (in Algorithm 2) to $T^{-1/4}$ (in Algorithm

3), which changes $T_s$ as well. The second is that the definition of $\hat{\theta}_s$ changes between the two algorithms. Note that the definition of $\hat{\theta}_s^{\mathrm{pre}}$ in Algorithm 3 is equivalent to the definition of $\hat{\theta}_s$ in Algorithm 2. This means that the definitions of $u_t^{\mathrm{safeU}}$ and $u_t^{\mathrm{safeL}}$ are the same in both Algorithm 2 and Algorithm 3. Of course, the changes in $\nu_T$ and the definition of $\hat{\theta}_s$ change the entire trajectory of the algorithm, which affects all of the other variables as well. However, all other differences in the algorithm trajectory can be derived from these two changes.

For the rest of Appendix F, let $C^{\mathrm{alg}}$ be the controller of Algorithm 3. Because Algorithm 3 and Algorithm 2 differ, we will now redefine the important events and lemmas from Appendix C with respect to Algorithm 3 (and the corresponding $\hat{\theta}_s$), and use this notation for the rest of Appendix F. Define $s_e = \log_2(T\nu_T^2) - 1$, and let

$$E_0 := \left\{ \forall s \le s_e : \|\theta^* - \hat{\theta}_s^{\mathrm{pre}}\|_\infty \le \epsilon_s \right\}. \tag{76}$$

By Lemma 23 we have that with probability $1 - o_T(1/T^2)$, $\|\theta^* - \hat{\theta}_s^{\mathrm{pre}}\|_\infty \le \epsilon_s$. Therefore,

$$\mathbb{P}(E_0) = 1 - o_T(1/T^2).$$

Also note that because $\|\hat{\theta}_s - \hat{\theta}_s^{\mathrm{pre}}\|_\infty \le \epsilon_s$ by construction, we have by the triangle inequality that under event $E_0$, $\|\theta^* - \hat{\theta}_s\|_\infty \le 2\epsilon_s$.

For the rest of this section, define

$$E_2 := E_0 \bigcap \left\{ \max_{s \in [0:s_e]} \epsilon_s = \tilde{O}_T(\nu_T) \right\}. \tag{77}$$

We also have the following equivalent result to Lemma 2, but with respect to the $\epsilon_s$ in Algorithm 3.

**Lemma 18.** *Under Assumptions 1–8, with probability $1 - o_T(1/T^2)$*

$$\max_{s \in [0:s_e]} \epsilon_s = \tilde{O}_T(\nu_T).$$

The proof of Lemma 2 relies only on the first $\nu_T$ steps and is written agnostic to the choice of $\nu_T$, and therefore the result of Lemma 18 follows directly from that proof. Lemma 18 implies that we have

$$\mathbb{P}(E_2) = 1 - o_T(1/T^2).$$

For this section, $E_1$ will still refer to the same event as in Equation (19). We also define the event $E_{\mathrm{safe}}$ the same way as in Equation (22) except with respect to the positions and controls of Algorithm 3, and finally we define the event $E = E_1 \cap E_2 \cap E_{\mathrm{safe}}$ (the same as in Appendix C.2). Therefore by a union bound we still have that $\mathbb{P}(E) = 1 - o_T(1/T^2)$. Using this new notation and Lemma 18, we can proceed to the main proof.

The safety of $C^{\mathrm{alg}}$ follows from an equivalent version of Lemma 1, except stated for Algorithm 3 instead of Algorithm 2. The proof follows as in the proof of Lemma 1 except using Lemma 18 instead of Lemma 2, and using the above definitions of $E_0, E_1$ and $E_2$ with respect to Algorithm 3. An equivalent statement of Lemma 3 holds except for the $u_t^{\mathrm{safeU}}$ and $u_t^{\mathrm{safeL}}$ coming from Algorithm 3. Note that the only place that the proof of Lemma 3

relies on $\nu_T$ is that it requires that $\epsilon_s = \tilde{O}_T(\nu_T)$ and that $\tilde{O}_T(\nu_T) = o_T(1/\log(T))$ at multiple points in the proof, which still holds under the new definitions of $E_2$ and $\nu_T$. Finally, as noted above, the $u_t^{\mathrm{safeU}}$ and $u_t^{\mathrm{safeL}}$ are constructed in the same way for both algorithms, and therefore the rest of the proof of Lemma 1 follows directly.

The rest of this section will focus on bounding the regret of Algorithm 3 to be $\tilde{O}_T(\sqrt{T})$ with probability $1 - o_T(1/T)$. Informally, the key idea behind the regret bound of Algorithm 3 is that with high probability, the uncertainty upper bound $\epsilon_s$ will decrease at a rate proportional to $1/\sqrt{T_s}$. This is formalized in Lemma 19.

**Lemma 19.** *Under Assumptions 1–9, given event $E$ with conditional probability $1 - o_T(1/T)$,*

$$\max_{s \in [0:s_e]} \epsilon_s \sqrt{T_s} = \tilde{O}_T(1).$$

The proof of Lemma 19 can be found in Appendix F.2.

Define event $E_3$ as the event

$$E_3 = \left\{ \max_{s \in [0:s_e]} \epsilon_s \sqrt{T_s} = \tilde{O}_T(1) \right\}.$$

By Lemma 19, $\mathbb{P}(E_3) = 1 - o_T(1/T)$. We can decompose the regret of Algorithm 3 into the same components of regret as in Appendix C.2. The first two propositions stated below are exactly equivalent to their counterparts in Appendix C.2.

**Proposition 7** (Regret from Warm-up Period). *Define $x'_0, x'_1, \ldots$ as the sequence of random variables that are the positions of the controller $C^{\mathrm{alg}}$ defined in Algorithm 3. Define $R_0$ as the cost of the first $1/\nu_T^2$ steps, i.e.*

$$R_0 = T \cdot J(\theta^*, C^{\mathrm{alg}}, T, 0, W) - \sum_{s=0}^{s_e} T_s \cdot J(\theta^*, C_s^{\mathrm{alg}}, T_s, x'_{T_s}, W_s). \tag{78}$$

*Then under Assumptions 1–8 and conditional on event $E$,*

$$R_0 \overset{a.s.}{\leq} \tilde{O}_T\left(\frac{1}{\nu_T^2}\right).$$

The proof of Proposition 7 can be found in Appendix F.3.

**Proposition 8** (Regret from Randomness). *Define $\hat{x}_{T_0}, \hat{x}_{T_0+1}, \ldots$ as the sequence of random variables representing the sequence of positions if the control at each time $t \geq T_0$ is $C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}(x_t)$ for $s = \lfloor \log_2(t\nu_T^2) \rfloor$ and starting at $\hat{x}_{T_0} = x'_{T_0}$. Define $R_2$ as*

$$R_2 := \sum_{s=0}^{s_e} T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, \hat{x}_{T_s}, W_s) - \sum_{s=0}^{s_e} \mathbb{E}\left[ T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \;\middle|\; \hat{\theta}_s \right].$$

*Then with conditional probability $1 - o_T(1/T)$ given event $E$,*

$$R_2 \leq \tilde{O}_T(\sqrt{T}). \tag{79}$$

The proof of Proposition 8 can be found in Appendix F.4. The next two propositions have different regret bounds than their counterparts in Appendix C.2.

**Proposition 9** (Regret from Non-optimal Controller with Sufficiently Large Noise). *Define $R_1$ as*

$$R_1 := \sum_{s=0}^{s_e} \mathbb{E}\left[ T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, 0, W_s) \,\Big|\, \hat{\theta}_s \right] - \mathbb{E}\left[ \sum_{s=0}^{s_e} T_s J(\theta^*, C^{\theta^*}_{K^*_s}, T_s, x^*_{T_s}, W_s) \right].$$

*Note that $W_s$ is independent of $\hat{\theta}_s$ by construction. Then under Assumptions 1–9 and conditional on event $E_2 \cap E_3$,*

$$R_1 \overset{a.s.}{\leq} \tilde{O}_T\left(\sqrt{T}\right). \tag{80}$$

The proof of Proposition 9 can be found in Appendix F.5.

**Proposition 10** (Regret from Enforcing Safety with Sufficiently Large Noise). *Define $\hat{x}_{T_0}, \hat{x}_{T_0+1}, \dots$ as the sequence of random variables representing the sequence of positions if the control at each time $t \geq T_0$ is $C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}(\hat{x}_t)$ for $s = \lfloor \log_2(t\nu_T^2) \rfloor$ and starting at $\hat{x}_{T_0} = x'_{T_0}$. Define $R_3$ as (the random variable)*

$$R_3 := \sum_{s=0}^{s_e} T_s J(\theta^*, C^{\mathrm{alg}}_s, T_s, x'_{T_s}, W_s) - \sum_{s=0}^{s_e} T_s J(\theta^*, C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}, T_s, \hat{x}_{T_s}, W_s).$$

*Then under Assumptions 1–9, with conditional probability $1 - o_T(1/T)$ given event $E \cap E_3$,*

$$R_3 \leq \tilde{O}_T(\sqrt{T}).$$

The proof of Proposition 10 can be found in Appendix F.6.

Using Equation (31) combined with Propositions 7, 8, 9 and 10, the total regret is upper bounded by the following conditioned on event $E \cap E_3$, with conditional probability $1 - o_T(1/T)$

$$T \cdot J(\theta^*, C^{\mathrm{alg}}, T) - T \cdot \bar{J}(\theta^*, C^{\theta^*}_{K_{\mathrm{opt}}(\theta^*, T)}, T) \leq R_0 + R_1 + R_2 + R_3 = \tilde{O}_T\left( \frac{1}{\nu_T^2} + \sqrt{T} \right).$$

Because $\nu_T = T^{-1/4}$ in Algorithm 3, this gives total regret of $\tilde{O}_T(\sqrt{T})$ conditional on $E_3 \cap E$. Since $\mathbb{P}(E_3) = 1 - o_T(1/T)$ and $\mathbb{P}(E) = 1 - o_T(1/T)$, a union bound gives that the regret of Algorithm 3 is $\tilde{O}_T(\sqrt{T})$ with unconditional probability $1 - o_T(1/T)$. $\qquad\square$

## F.2  Proof of Lemma 19(Uncertainty bounds using boundary times)

*proof.* To prove this lemma, we will show that the controller $C^{\mathrm{alg}}$ uses the control $u_i^{\mathrm{safeU}}$ "sufficiently frequently". Let $S_t$ be the set of times $i < t$ when the control used by Algorithm 3 is $u_i^{\mathrm{safeU}}$. Formally, if $u'_0, u'_1, \dots u'_{T-1}$ are the sequence of controls used by $C^{\mathrm{alg}}$, then

$$S_t = \{i < t : u'_i = u_i^{\mathrm{safeU}}\}. \tag{81}$$

**Lemma 20.** *Under Assumptions 1–9 and conditional on event $E$ with conditional probability* $1 - o_T(1/T)$,

$$\min_{s \in [1:s_e]} \frac{|S_{T_s}|}{T_s} = \Omega_T(1).$$

The proof of Lemma 20 can be found in Appendix F.8. Equipped with the fact that $|S_t|$ scales linearly with $t$ from Lemma 20, we need the following result that will lower the upper bound for $\epsilon_s$.

**Lemma 21.** *Under Assumptions 1–9 and conditional on event $E$ with conditional probability* $1 - o_T(1/T)$,

$$\max_{s \in [0:s_e]} \epsilon_s \sqrt{|S_{T_s}|} = \tilde{O}_T(1).$$

The proof of Lemma 21 can be found in Appendix G.3. To see that $\epsilon_0 \sqrt{T_0} = \tilde{O}_T(1)$, note that $\sqrt{T_0} = 1/\nu_T$ and Lemma 18 imply that conditional on event $E$, $\epsilon_0 = \tilde{O}_T(\nu_T)$. For $s > 0$, a union bound combining Lemma 20 with Lemma 21 gives the desired result that conditioned on event $E$ with conditional probability $1 - o_T(1/T)$,

$$\max_{s \in [0:s_e]} \epsilon_s \sqrt{T_s} = \tilde{O}_T(1).$$

$\square$

## F.3    Proof of Proposition 7

*proof.* The proof of Proposition 7 follows the same as the proof of Proposition 3. The proof of Proposition 3 relies on the fact that the controller is safe for dynamics $\theta^*$ conditional on event $E$. This is still true by construction of event $E$, and therefore the result follows directly. $\square$

## F.4    Proof of Proposition 8

*proof.* Note that this statement is exactly the same as the statement of Proposition 5 except for Algorithm 3. The proof of Proposition 5 relies on Lemmas 6 and 7. Define the event $E_2^s$ as

$$E_2^s = \left\{ \|\hat{\theta}_s - \theta^*\|_\infty \leq 2 \cdot \epsilon_s \leq 2 c_T \cdot \nu_T \right\}, \tag{82}$$

where the $c_T = \tilde{O}_T(1)$ from Lemma 18. Note that we still have $\mathbb{P}(E_2^s) \geq \mathbb{P}(E_2) \geq 1 - o_T(1/T^2)$. An analogous version of Lemma 6 holds with this new definition of $E_2^s$ for Algorithm 3. Examining Lemma 6, the proof relies on $\hat{\theta}_s$ and $\nu_T$ through Lemma 15. A version of Lemma 15 holds with the exact same statement with the new definition of $E_2^s$. Examining the proof of Lemma 15, we must have that under event $E_2^s$, $\|\hat{\theta}_s - \theta^*\|_\infty \leq \tilde{O}_T(\nu_T) \leq \min(\epsilon_{A8}, \frac{1}{\log(T)})$ in order to apply Lemmas 8 and 4, and this holds for $\nu_T = T^{-1/4}$. Therefore, we have shown the equivalent version of Lemma 6 for Algorithm 3.

Similarly, an analogous version of Lemma 7 holds for Algorithm 3. Lemma 7 depends on $\hat{\theta}_s$ and $\nu_T$ only in that it uses $\|\theta^* - \hat{\theta}_s\|_\infty \leq 1/\log(T)$ conditional on event $E_2^s$, which still holds by construction of $E_2^s$ for $\nu_T = T^{-1/4}$ and sufficiently large $T$.

Now that we have shown that equivalent versions of Lemmas 6 and 7 still hold, we can return to the proof of Proposition 5. Outside of the two lemmas discussed above, the only places in the proof that depend on the choice of $\nu_T$ and $\hat{\theta}_s$ is that $s_e = \tilde{O}_T(1)$ is still true in Equation (44) and that conditional on event $E$, $\|\hat{\theta}_s - \theta^*\|_\infty \leq \tilde{O}_T(\nu_T) \leq \min(\epsilon_{A8}, \frac{1}{\log(T)})$ in order to apply Lemmas 4 and 8. As both of these still hold for the new definition of $E$ and for $\nu_T = T^{-1/4}$, we are done. □

## F.5 Proof of Proposition 9

*proof.* The proof of Proposition 9 will mostly follow as in the proof of Proposition 4. The proof of Proposition 4 relies on Lemma 5. An equivalent version of Lemma 5 holds for Algorithm 3, where the only difference is that the $T_s$ are now defined differently. To see this, note that the proof of Lemma 5 works for any $T_s \leq T$, and therefore the proof follows exactly the same.

Returning to the proof of Proposition 4, we can still apply Assumption 7 under the event $E_2^s$ as defined in Equation (82). Looking at the last block of equations in Proposition 4, we can follow the logic exactly and pick up from the second to last line. Applying Lemma 19, conditional on $E \cap E_3$,

$$\sum_{s=0}^{s_e} \mathbb{E}\left[T_s J(\theta^*, C_{K_{\text{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}, T_s, 0, W_s) \,\Big|\, \hat{\theta}_s\right] - \mathbb{E}\left[\sum_{s=0}^{s_e} T_s J(\theta^*, C_{K_s^*}^{\theta^*}, T_s, x_{T_s}^*, W_s)\right]$$

$$\leq \tilde{O}_T(1) + \tilde{O}_T\left(\sum_{s=0}^{s_e} T_s \epsilon_s + \frac{T_s}{T^2}\right)$$

$$= \tilde{O}_T(1) + \tilde{O}_T\left(\sum_{s=0}^{s_e} T_s \tilde{O}_T\left(\frac{1}{\sqrt{T_s}}\right) + \frac{T_s}{T^2}\right) \qquad \text{Lemma 19}$$

$$= \tilde{O}_T(\sqrt{T}).$$

□

## F.6 Proof of Proposition 10

*proof.* The proof of Proposition 10 will mostly follow as in the proof of Proposition 6. The proof of Proposition 6 relies on Lemmas 9 and 10. We will show that equivalent versions of these lemmas hold for Algorithm 3.

Starting with Lemma 9, an equivalent version holds for the $u_t^{\text{safeU}}$ and $u_t^{\text{safeL}}$ defined in Algorithm 3 and $C^{\text{alg}}$ as the controller of Algorithm 3. Looking at the proof of Lemma 9, the main tool is Lemma 16. An equivalent version of Lemma 16 holds for Algorithm 3. Looking at the proof of Lemma 16, the dependency on $\hat{\theta}_s$ and $\nu_T$ is that we must have that conditional on event $E$, $\|\hat{\theta}_s - \theta^*\|_\infty \leq \tilde{O}_T(\nu_T) \leq \min(\epsilon_{A8}, \frac{1}{\log(T)})$ in order to apply Lemmas 4 and 8. The union bound at the end of the proof also relies on $s_e = \tilde{O}_T(1)$, which also does still hold. Returning to the proof of the equivalent of Lemma 9 for Algorithm 3, we again need that conditional on event $E$, $\|\hat{\theta}_s - \theta^*\|_\infty \leq \tilde{O}_T(\nu_T) \leq \min(\epsilon_{A8}, \frac{1}{\log(T)})$ in order to apply

Lemmas 4 and 8. Once again using that $s_e = \tilde{O}_T(1)$, the rest of the proof of Lemma 9 can be directly applied.

An equivalent version of Lemma 10 also holds when $C^{\mathrm{alg}}$ is the controller of Algorithm 3 with $\nu_T = T^{-1/4}$. We defer the proof of this to Appendix F.7.

Now we can return to the proof of Proposition 6 and show that a slight modification gives the desired result. Looking at the last set of equations, we can pick up from the third line and apply Lemma 19 to get that, conditional on event $E \cap E_3$,

$$
\begin{aligned}
R_3 &\leq \tilde{O}_T \left( \sum_{s=0}^{s_e} T_s \epsilon_s \right) + \sum_{s=0}^{s_e} \sum_{t=T_s}^{T_{s+1}-1} X_t^U \cdot \tilde{O}_T(\epsilon_s) + X_t^L \cdot \tilde{O}_T(\epsilon_s) \\
&\leq \tilde{O}_T \left( \sum_{s=0}^{s_e} T_s \epsilon_s \right) \\
&\leq \sum_{s=0}^{s_e} T_s \cdot \tilde{O}_T \left( \frac{1}{\sqrt{T_s}} \right) \qquad\qquad\qquad\qquad \text{Lemma 19} \\
&= \tilde{O}_T(\sqrt{T}).
\end{aligned}
$$

The last line follows from the fact that for all $s$, $T_s \leq T$ and that $s_e = \tilde{O}_T(1)$. $\qquad\square$

## F.7 Proof of Equivalent Version of Lemma 10 for Algorithm 3

Examining the proof of Lemma 10, the main change when using Algorithm 3 is that we now have that under event $E$ and for sufficiently large $T$, $\|\theta^* - \hat{\theta}_s\|_\infty \leq 2\epsilon_s$ (while for Algorithm 2 there was no factor of 2). Because $\nu_T = T^{-1/4}$, this still allows us to apply Lemma 4. Picking up the proof of Lemma 10 directly before Equation (61), the extra factor of 2 mentioned above will result in the following changes.

By the construction of Algorithm 3, $u_t^{\mathrm{safeU}}$ satisfies, for some $\theta$ such that $\|\theta - \hat{\theta}_s^{\mathrm{pre}}\|_\infty \leq \epsilon_s$,

$$ ax_t' + bu_t^{\mathrm{safeU}} = D_{\mathrm{U}}. \tag{83} $$

Under event $E$, $\|\theta^* - \hat{\theta}_s\|_\infty \leq 2\epsilon_s$ and $\|\hat{\theta}_s - \hat{\theta}_s^{\mathrm{pre}}\|_\infty \leq \epsilon_s$, which implies that $\|\theta^* - \theta\|_\infty \leq 4\epsilon_s \leq \tilde{O}_T(\nu_T) \leq 1/\log(T)$ for sufficiently large $T$. Therefore, applying Lemma 4 gives that under event $E$ and for sufficiently large $T$,

$$
\begin{aligned}
D_{\mathrm{U}} &\geq a^* x_t' + b^* u_t^{\mathrm{safeU}} & u_t^{\mathrm{safeU}} \text{ safe for } \theta^* \\
&\geq ax_t' + bu_t^{\mathrm{safeU}} - |u_t^{\mathrm{safeU}}|4\epsilon_s - |x_t'|4\epsilon_s & \|\theta^* - \theta\|_\infty \leq 4\epsilon_s \\
&\geq D_{\mathrm{U}} - 8B_x \epsilon_s. & \text{Equations (58),(59), and (83)} \tag{84}
\end{aligned}
$$

If $u_t^{\mathrm{safeU}} \leq C_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}(x_t')$, then there must exist some $\theta$ such that $\|\hat{\theta}_s^{\mathrm{pre}} - \theta\|_\infty \leq \epsilon_s$ and

$$ ax_t' + bC_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s)}^{\hat{\theta}_s}(x_t') \geq D_{\mathrm{U}}. \tag{85} $$

By the same logic as above, under event $E$, $\|\theta^* - \theta\|_\infty \le 4\epsilon_s \le \tilde{O}_T(\nu_T) \le 1/\log(T)$ for sufficiently large $T$, therefore under event $E$ and for sufficiently large $T$,

$$a^* x'_t + b^* C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x'_t)$$
$$\ge a x'_t + b C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x'_t) - 4\epsilon_s |x'_t| - 4\epsilon_s \left| C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x'_t) \right|$$
$$\ge D_{\text{U}} - 8 B_x \epsilon_s. \qquad \text{Eqs (58),(60), and (85)} \qquad (86)$$

Finally, because $C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x'_t)$ is safe for dynamics $\hat{\theta}_s$,

$$\hat{a}_s x'_t + \hat{b}_s C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x'_t) \le D_{\text{U}}. \qquad (87)$$

Using that under event $E$, $\|\theta^* - \hat{\theta}_s\|_\infty \le 2\epsilon_s \le \tilde{O}_T(\nu_T) \le 1/\log(T)$ for sufficiently large $T$, Equations (58), (60), and (87) imply that under event $E$ and for sufficiently large $T$,

$$a^* x'_t + b^* C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x'_t) \le D_{\text{U}} + 4 B_x \epsilon_s. \qquad (88)$$

Combining Equations (86) and (88), if $u_t^{\text{safeU}} \le C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x'_t)$ then under event $E$ and for sufficiently large $T$,

$$D_{\text{U}} - 8 B_x \epsilon_s \le a^* x'_t + b^* C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x'_t) \le D_{\text{U}} + 4 B_x \epsilon_s.$$

Combining this with Equation (84) gives that under event $E$ and for sufficiently large $T$,

$$|(a^* x'_t + b^* u_t^{\text{safeU}}) - (a^* x'_t + b^* C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x'_t))| \le 12 B_x \epsilon_s.$$

This implies the desired result that under event $E$ and for sufficiently large $T$,

$$|u_t^{\text{safeU}} - C^{\hat{\theta}_s}_{K_{\text{opt}}(\hat{\theta}_s, T_s)}(x'_t)| \le 12 B_x \epsilon_s / b^*.$$

## F.8  Proof of Lemma 20

*proof.* In this proof, we will use the following result about the times the algorithm chooses control $u_t^{\text{safeU}}$.

**Lemma 22.** *Let $x'_t$, $u'_t$ be the positions and controls of controller $C^{\text{alg}}$ at time $t$. For $t \ge T_0$, let $s_t = \lfloor \log_2(t \nu_T) \rfloor$. Then under Assumptions 1–9, there exists a $P_{s_t}(\hat{\theta}_{s_t}, \epsilon_{s_t})$ such that*

$$\{x'_t \ge P_{s_t}(\hat{\theta}_{s_t}, \epsilon_{s_t})\} \subseteq \{u'_t = u_t^{\text{safeU}}\},$$

*and such that conditional on event $E$, we have $P_{s_t}(\hat{\theta}_{s_t}, \epsilon_{s_t}) \le P(\theta^*, K_{\text{opt}}(\theta^*, T_{s_t}), D_{\text{U}})$.*

The proof of Lemma 22 can be found in Appendix F.9. Recall that $\{i \in S_{T_s}\} = \{u'_i = u_i^{\text{safeU}}\}$. Therefore, for $i \in [T_s : T_{s+1} - 1]$, Lemma 22 implies that

$$\{x'_i \ge P(\theta^*, K_{\text{opt}}(\theta^*, T_s), D_{\text{U}})\} \cap E \subseteq \{x'_i \ge P_{s_t}(\hat{\theta}_{s_t}, \epsilon_{s_t})\} \cap E$$
$$\subseteq \{u'_i = u_i^{\text{safeU}}\} \cap E$$
$$= \{i \in S_{T_s}\} \cap E. \qquad (89)$$

By Assumption 9, for any $x'_{i-1}, u'_{i-1}$ satisfying $a^* x'_{i-1} + b^* u'_{i-1} \in [D_L, D_U]$, we have that

$$
\begin{aligned}
&\mathbb{P}\left(x'_i \geq P(\theta^*, K_{opt}(\theta^*, T_{s_t}), D_U) \mid x'_{i-1}, u'_{i-1}\right) \\
&\geq \mathbb{P}\left(w_i \geq P(\theta^*, K_{opt}(\theta^*, T_{s_t}), D_U) - D_L\right) \\
&\geq \epsilon_{A9}. \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{Assumption 9} \qquad (90)
\end{aligned}
$$

Because $\mathbb{P}(E) \geq 1 - o_T(1/T)$, this implies for sufficiently large $T$ and for any $x'_{i-1}, u'_{i-1}$ satisfying $a^* x'_{i-1} + b^* u'_{i-1} \in [D_L, D_U]$,

$$
\begin{aligned}
\mathbb{P}\left(x'_i \geq P(\theta^*, K_{opt}(\theta^*, T_{s_t}), D_U) \mid x'_{i-1}, u'_{i-1}, E\right) &\geq \epsilon_{A9} - \mathbb{P}(\neg E) \\
&\geq \epsilon_{A9} - o_T(1/T) \\
&\geq \frac{\epsilon_{A9}}{2}. \qquad\qquad (91)
\end{aligned}
$$

Also, recall that conditional on event $E$, $C^{alg}$ is safe for dynamics $\theta^*$ for all $T$ steps, therefore conditional on event $E$, for all $i \geq 1$, $D_L \leq a^* x'_{i-1} + b^* u'_{i-1} \leq D_U$. Therefore, for $T_1 \leq i < T_s$ and sufficiently large $T$,

$$
\begin{aligned}
&\mathbb{P}\left(i \in S_{T_s} \mid x'_0, x'_1, ..., x'_{i-1}, u'_0, u'_1, ..., u'_{i-1}, E\right) \\
&\geq \mathbb{P}\left(x'_i \geq P(\theta^*, K_{opt}(\theta^*, T_s), D_U) \mid x'_0, x'_1, ..., x'_{i-1}, u'_0, u'_1, ..., u'_{i-1}, E\right) \quad \text{Equation (89)} \\
&\geq \mathbb{P}\left(x'_i \geq P(\theta^*, K_{opt}(\theta^*, T_s), D_U) \mid x'_{i-1}, u'_{i-1}, E\right) \\
&\geq \frac{\epsilon_{A9}}{2}. \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{Equation (91)} \quad (92)
\end{aligned}
$$

Defining $X_i = 1_{i \in S_{T_s}}$, the above equation is equivalent to

$$
\mathbb{E}\left[X_i \mid x'_0, x'_1, ..., x'_{i-1}, u'_0, u'_1, ..., u'_{i-1}, E\right] \geq \frac{\epsilon_{A9}}{2}.
$$

Therefore, we can conclude that conditional on event $E$, $\sum_{i=T_0}^{T_s-1} X_i$ is stochastically dominated by $\sum_{i=T_0}^{T_s-1} Y_i$, where $Y_i$ are i.i.d. Bernoulli random variables that are equal to 1 with probability $\epsilon_{A9}/2$. By this coupling argument and Hoeffding's inequality, for $s \geq 1$, conditional on event $E$ with conditional probability $1 - o_T(1/T^2)$,

$$
|S_{T_s}| = \sum_{i=T_0}^{T_s-1} X_i \geq \sum_{i=T_0}^{T_s-1} Y_i \geq \frac{\epsilon_{A9}}{2}(T_s - T_0) - \log(T)\sqrt{T_s - T_0} \geq \frac{\epsilon_{A9}}{4} \cdot (T_s - T_0) \geq \frac{\epsilon_{A9}}{8} \cdot T_s, \quad (93)
$$

where the second to last inequality comes from for sufficiently large $T$ and $s \geq 1$, $T_s - T_0 \geq \sqrt{T}$ and therefore $\sqrt{T_s - T_0} \geq \frac{4\log(T)}{\epsilon_{A9}}$. The last inequality comes from the fact that $T_s - T_0 \geq \frac{T_s}{2}$ by the definition of $T_s$ for $s \geq 1$. A union bound over all $s \in [1 : s_e]$ gives that conditional on event $E$ with conditional probability $1 - o_T(1/T)$,

$$
\min_{s \in [1:s_e]} \frac{|S_{T_s}|}{T_s} \geq \frac{\epsilon_{A9}}{8}.
$$

$\square$

## F.9 Proof of Lemma 22

*proof.* Defining $P_s(\hat{\theta}_s, \epsilon_s)$ as

$$P_s(\hat{\theta}_s, \epsilon_s) = \min_{\|\theta - \hat{\theta}_s^{\mathrm{pre}}\|_\infty \leq \epsilon_s} P(\theta, K_{\mathrm{opt}}(\hat{\theta}_s, T_s), D_{\mathrm{U}}),$$

we have by definition of $u_t^{\mathrm{safeU}}$ in Algorithm (3) that

$$\{x_t' \geq P_s(\hat{\theta}_s, \epsilon_s)\} = \left\{x_t' \geq \min_{\|\theta - \hat{\theta}_s^{\mathrm{pre}}\|_\infty \leq \epsilon_s} P(\theta, K_{\mathrm{opt}}(\hat{\theta}_s, T_s), D_{\mathrm{U}})\right\} \subseteq \{u_t' = u_t^{\mathrm{safeU}}\}. \tag{94}$$

Under event $E$, $\|\theta^* - \hat{\theta}_s^{\mathrm{pre}}\|_\infty \leq \epsilon_s$, therefore

$$\min_{\|\theta - \hat{\theta}_s^{\mathrm{pre}}\|_\infty \leq \epsilon_s} P(\theta, K_{\mathrm{opt}}(\theta^*, T_s), D_{\mathrm{U}}) \leq P(\theta^*, K_{\mathrm{opt}}(\theta^*, T_s), D_{\mathrm{U}}). \tag{95}$$

Therefore, we can conclude that conditional on event $E$,

$$
\begin{aligned}
P_s(\hat{\theta}_s, \epsilon_s) &= \min_{\|\theta - \hat{\theta}_s^{\mathrm{pre}}\|_\infty \leq \epsilon_s} P(\theta, K_{\mathrm{opt}}(\hat{\theta}_s, T_s), D_{\mathrm{U}}) \\
&\leq \min_{\|\theta - \hat{\theta}_s^{\mathrm{pre}}\|_\infty \leq \epsilon_s} P(\theta, K_{\mathrm{opt}}(\theta^*, T_s), D_{\mathrm{U}}) && \text{Choice of } \hat{\theta}_s \\
&\leq P(\theta^*, K_{\mathrm{opt}}(\theta^*, T_s), D_{\mathrm{U}}). && \text{Equation (95)}
\end{aligned}
$$

$\square$

# G   Uncertainty Bounds

## G.1   Tools for Uncertainty Bounds

The proofs of uncertainty bounds will rely on the following result from Abbasi-Yadkori and Szepesvári [2011].

**Lemma 23** (Derived from Theorem 1 in Abbasi-Yadkori and Szepesvári [2011]). *Suppose $x_t$ and $u_t$ are respectively the position and control at time $t$ when using an arbitrary controller $C$ starting at position $x_0 = 0$. Define $z_t = (x_t, u_t)$ and let $\lambda > 0$. Let $Z_t \in \mathbb{R}^{t \times 2}$ where the ith row is $z_{i-1}$, let $X_t \in \mathbb{R}^{t \times 1}$ where the ith element is $x_i$, and let $I \in \mathbb{R}^{2 \times 2}$ be the identity matrix. Then under Assumptions 1–3, with probability $1 - o_T\left(\frac{1}{T^2}\right)$ the following holds for all $1 \leq t \leq T - 1$ and for any $S \subseteq [0 : t - 1]$:*

$$\|\theta^* - (Z_t^\top Z_t + \lambda I)^{-1} Z_t^\top X_t\|_\infty \leq \sqrt{\frac{\max((V_t^S)_{11}, (V_t^S)_{22})}{\det(V_t^S)}} B_t, \tag{96}$$

*where $V_t^S = \lambda I + \sum_{s=0}^{t-1} z_s z_s^\top 1_{s \in S}$, $B_t = \alpha\sqrt{\log\left(\det\left(V_t^{[0:t-1]}\right)\right) + \log(\lambda^2) + 2\log(T^2)} + \sqrt{\lambda}(\bar{a}^2 + \bar{b}^2)$, and $\alpha$ is from the subgaussian assumption on the noise distribution $\mathcal{D}$, which implies that there exists an $\alpha$ such that $\mathbb{E}_{w \sim \mathcal{D}}[\exp(\gamma w)] \leq \exp(\gamma^2 \alpha^2 / 2)$ for any $\gamma \in \mathbb{R}$.*

Lemma 23 can be directly derived from Theorem 1 in Abbasi-Yadkori and Szepesvári [2011] as shown in Appendix G.5. The other tool that will be shared by the proofs in the following sections is the following anti-concentration inequality of the sum of non-negative random variables.

**Lemma 24.** *For $p \in (0, 1]$ and $1 \leq n \leq T$, suppose $X_0, ..., X_{n-1}$ are non-negative random variables such that $(X_0, ..., X_{i-1})$ is a deterministic function of the random variable set $F_i$ for all $i \in [1 : n]$ and $F_i \subseteq F_{i+1}$. Let the set $S_n \subseteq [0 : n-1]$ be a random variable such that the event $\{i \in S_n\}$ is a deterministic function of $F_{i+1}$. For $i \in [0 : n-1]$, define $S_i = \{k < i : k \in S_n\}$, therefore $S_i$ is a deterministic function of $F_i$. Let $E^*$ be an event such that for all $i \in [0 : n-1]$,*

$$\mathbb{E}\left[X_i \mid F_i, E^*, i \in S_n\right] \geq c \cdot |S_i|, \tag{97}$$

*where $c > 0$ is non-random. Furthermore, assume that conditional on $E^*$:*

$$0 \overset{a.s.}{\leq} X_i \overset{a.s.}{\leq} \frac{c|S_i|}{2p}. \tag{98}$$

*Then conditional on event $E^*$, with conditional probability $1 - o_T(1/T^2)$,*

$$\sum_{i=0}^{n-1} X_i \geq \frac{c}{4} \left(\max(\lfloor p|S_n| - \log(T)\sqrt{|S_n|}\rfloor, 1)\right) \left(\max(\lfloor p|S_n| - \log(T)\sqrt{|S_n|}\rfloor, 1) - 1\right).$$

The proof of Lemma 24 can be found in Appendix G.6.

## G.2 Proof of Lemma 2

*proof.* For the rest of this proof, $x_t$ and $u_t$ are respectively the position and control at time $t$ of controller $C^{\text{alg}}$ that corresponds to Algorithm 2 starting at $x_0 = 0$. Recall that Lemma 2 was stated and used in Appendix C with respect to Algorithm 2, therefore all events and variables in this subsection refer to those defined with respect to Algorithm 2. To prove Lemma 2, we will use Lemma 23 applied to $S = [0 : \frac{1}{\nu_T^2} - 1]$. The goal will be to bound the right side of Equation (96) for this choice of $S$. Consider any fixed arbitrary $s \in [0 : s_e]$ and the corresponding matrix $V_{T_s}^S$. Define $N$ as the event that for all $i < 1/\nu_T^2$, the control $u_i$ is safe for dynamics $\theta^*$. Note that we showed in Lemma 1 that $\mathbb{P}(N) \geq \mathbb{P}(E_{\text{safe}}) = 1 - o_T(\frac{1}{T^2})$. Under event $N \cap E_1$, we can apply Lemma 4 to get the following equations for sufficiently large $T$:

$$(V_{T_s}^S)_{11} = \lambda + \sum_{i=0}^{\frac{1}{\nu_T^2} - 1} x_i^2 \leq \lambda + \frac{1}{\nu_T^2} B_x^2 \tag{99}$$

$$(V_{T_s}^S)_{22} = \lambda + \sum_{i=0}^{\frac{1}{\nu_T^2} - 1} u_i^2 \leq \lambda + \frac{1}{\nu_T^2} B_x^2 \tag{100}$$

$$(V_{T_s}^S)_{12}^2 = \left( \sum_{i=0}^{\frac{1}{\nu_T^2}-1} u_i x_i \right)^2. \tag{101}$$

We can now compute $(V_{T_s}^S)_{22}(V_{T_s}^S)_{11} - (V_{T_s}^S)_{12}^2$. Recall that for the first $1/\nu_T^2$ steps of Algorithm 2, the control is $u_i = C^{\text{init}}(x_i) + \frac{\phi_i}{\log(T)}$ where $\phi_i$ is i.i.d. from the Rademacher distribution and independent from the noise random variables.

$$(V_{T_s}^S)_{22}(V_{T_s}^S)_{11} - (V_{T_s}^S)_{12}^2$$

$$= \left( \lambda + \sum_{i=0}^{\frac{1}{\nu_T^2}-1} u_i^2 \right)\left( \lambda + \sum_{i=0}^{\frac{1}{\nu_T^2}-1} x_i^2 \right) - \left( \sum_{i=0}^{\frac{1}{\nu_T^2}-1} u_i x_i \right)^2 \qquad \text{Equations (99) (100) (101)}$$

$$\geq \left( \sum_{i=0}^{\frac{1}{\nu_T^2}-1} u_i^2 \right)\left( \sum_{i=0}^{\frac{1}{\nu_T^2}-1} x_i^2 \right) - \left( \sum_{i=0}^{\frac{1}{\nu_T^2}-1} u_i x_i \right)^2$$

$$= \sum_{i<j}^{\frac{1}{\nu_T^2}-1} (u_i x_j - u_j x_i)^2$$

$$= \sum_{i<j}^{\frac{1}{\nu_T^2}-1} \left( u_i x_j - C^{\text{init}}(x_j) x_i + \frac{\phi_j}{\log(T)} x_i \right)^2. \tag{102}$$

Conditional on $N \cap E_1$, for all $i < 1/\nu_T^2$, we have $|u_i|, |x_i| \leq B_x$ by Lemma 4. Define the random variable $X_j$ as

$$X_j = \sum_{i=0}^{j-1} (u_i x_j - u_j x_i)^2$$

$$= \sum_{i=0}^{j-1} \left( u_i x_j - C^{\text{init}}(x_j) x_i + \left( \frac{\phi_j}{\log(T)} \right) x_i \right)^2$$

$$\leq 4j B_x^4. \qquad \text{Conditional on } N \cap E_1 \text{ by Lemma 4} \tag{103}$$

We will use the following lemma to lower bound the conditional expectation of $X_j$.

**Lemma 25.** *Under Assumptions 1–3, let $x_0, x_1, ..., x_T$ be the positions of the controller $C^{\text{alg}}$ starting at $x_0 = 0$. Then there exists an event $E_{L25}$ such that $\mathbb{P}(E_{L25}) = 1 - o_T(1/T^2)$ and for sufficiently large $T$ conditional on $E_{L25}$, for all $j \geq \log^8(T)$,*

$$\sum_{i=0}^{j-1} x_i^2 \geq \frac{j}{2\log^2(T)}. \tag{104}$$

The proof of Lemma 25 can be found in Appendix G.7. Now define $E^* = N \cap E_1 \cap E_{L25}$. Note that $\mathbb{P}(E^*) = 1 - o_T(1/T^2)$ by a union bound. Because $\phi_j$ is a Rademacher random variable, we therefore have that $\mathbb{P}(\phi_j = 1 \mid E^*) = 1/2 - o_T(1/T^2)$ and $\mathbb{P}(\phi_j = -1 \mid E^*) = 1/2 - o_T(1/T^2)$. This implies that $|\mathbb{E}[\phi_j \mid E^*]| = o_T(1/T^2)$, and therefore for sufficiently large $T$, we have $\mathrm{Var}\,[\phi_j \mid E^*] \geq 1/2$. Then we can bound the conditional expectation of $X_j$ under event $E^*$ as follows for all $j \geq \log^8(T)$ and for sufficiently large $T$. Define $F_j = \{x_0, u_0, ..., x_{j-1}, u_{j-1}, x_j\}$. Then we have

$$
\begin{aligned}
\mathbb{E}[X_j \mid F_j, E^*] &= \sum_{i=0}^{j-1} \mathbb{E}\left[\left(u_i x_j - C^{\mathrm{init}}(x_j)x_i + \left(\frac{\phi_j}{\log(T)}\right)x_i\right)^2 \middle| F_j, E^*\right] \\
&\geq \sum_{i=0}^{j-1} \mathrm{Var}\left[u_i x_j - C^{\mathrm{init}}(x_j)x_i + \left(\frac{\phi_j}{\log(T)}\right)x_i \middle| F_j, E^*\right] \\
&= \sum_{i=0}^{j-1} x_i^2 \,\mathrm{Var}\left[\frac{\phi_j}{\log(T)} \middle| F_j, E^*\right] \\
&= \sum_{i=0}^{j-1} x_i^2 \,\mathrm{Var}\left[\frac{\phi_j}{\log(T)} \middle| E^*\right] && \phi_j \text{ is ind. of } F_j \\
&= \frac{\mathrm{Var}\,[\phi_j \mid E^*]}{\log^2(T)} \sum_{i=0}^{j-1} x_i^2 \\
&\geq \frac{1}{2\log^2(T)} \sum_{i=0}^{j-1} x_i^2 \\
&\geq \frac{j}{4\log^4(T)}. && E^* \subseteq E_{L25} \quad (105)
\end{aligned}
$$

Therefore, we can apply Lemma 24 to $X_{\log^8(T)}, X_{\log^8(T)+1}, ..., X_{1/\nu_T^2 - 1}$ with $n = 1/\nu_T^2 - \log^8(T)$, $p = \frac{1}{32 B_x^4 \log^4(T)}$, $F_i = \{x_0, u_0, ..., u_{i-1}, x_i\}$, $S_n = [0 : n-1]$, and $c = \frac{1}{4\log^4(T)}$. Note that this choice of $p$ is less than 1 for sufficiently large $T$.

We will also use that for sufficiently large $T$, $n = 1/\nu_T^2 - \log^8(T) = T^{2/3} - \log^8(T) \geq 4\log^2(T)/p^2$. This implies that for sufficiently large $T$,

$$
pn - \log(T)\sqrt{n} \geq pn/2 = \frac{1/\nu_T^2 - \log^8(T)}{64 B_x^4 \log^4(T)} \geq 1. \tag{106}
$$

Recall by Equation (103) that under event $E^*$, the $X_j$ are bounded by $0 \leq X_j \leq 4jB_x^4 = \frac{c}{2p} \cdot j$. Lemma 24 gives that for sufficiently large $T$ and conditional on event $E^*$ with condi-

tional probability $1 - o_T(1/T^2)$,

$$(V^S_{T_s})_{22}(V^S_{T_s})_{11} - (V^S_{T_s})^2_{12}$$

$$\geq \sum_{j=0}^{\frac{1}{\nu^2_T}-1} X_j \qquad\qquad\qquad \text{Equation (102)}$$

$$\geq \sum_{j=\log^8(T)}^{\frac{1}{\nu^2_T}-1} X_j \qquad\qquad\qquad X_i \geq 0$$

$$\geq \frac{1}{16\log^4(T)}\left(\max(\lfloor pn - \log(T)\sqrt{n}\rfloor, 1) - 1\right)\left(\max(\lfloor pn - \log(T)\sqrt{n}\rfloor, 1)\right) \quad \text{Lemma 24}$$

$$\geq \frac{1}{16\log^4(T)}\left(\left\lfloor\frac{1/\nu^2_T - \log^8(T)}{64B^4_x\log^4(T)}\right\rfloor - 1\right)\left(\left\lfloor\frac{1/\nu^2_T - \log^8(T)}{64B^4_x\log^4(T)}\right\rfloor\right) \qquad \text{Equation (106)}$$

$$= \Omega_T\left(\frac{\frac{1}{\nu^4_T}}{B^8_x\log^{12}(T)}\right). \qquad\qquad\qquad (107)$$

Finally, we need to bound the quantity $B_{T_s}$ from Lemma 23. The only non-constant term in $B_{T_s}$ is $\sqrt{\log(\det(\lambda I + \sum_{i=0}^{T_s-1} z_i z_i^\top)) + 2\log(T^2)}$. Define $V_{T_s} = \lambda I + \sum_{i=0}^{T_s-1} z_i z_i^\top$. Conditional on event $N \cap E_1$, we have by Lemma 4 that $(V_{T_s})_{22} \leq \lambda + TB^2_x$ and $(V_{T_s})_{11} \leq \lambda + TB^2_x$. Therefore, conditional on event $N \cap E_1$,

$$\sqrt{\log\left(\det\left(\lambda I + \sum_{i=0}^{T_s-1} z_i z_i^\top\right)\right) + 2\log(T^2)} \leq \sqrt{\log((V_{T_s})_{11}(V_{T_s})_{22}) + 2\log(T^2)}$$

$$\leq \sqrt{\log\left((\lambda + TB^2_x)^2\right) + 2\log(T^2)}$$

$$= \tilde{O}_T(1). \qquad\qquad\qquad (108)$$

Now, combining Lemma 23 and Equations (99), (100), (107), and (108) gives that conditional on event $E^*$ with conditional probability $1 - o_T(1/T^2)$, for all $s \in [0 : s_e]$,

$$\epsilon^2_s \leq \frac{\max((V^S_{T_s})_{11}, (V^S_{T_s})_{22})}{\det(V^S_{T_s})}B^2_{T_s} = \frac{\left(\lambda + \left(\frac{1}{\nu^2_T}\right)B^2_x\right)\tilde{O}_T(1)}{\Omega_T\left(\frac{\left(\frac{1}{\nu^2_T}\right)^2}{B^8_x\log^{12}(T)}\right)} = \tilde{O}_T\left(\nu^2_T\right).$$

Because $\mathbb{P}(E^*) = 1 - o_T(1/T^2)$, this gives the desired result with unconditional probability $1 - o_T(1/T^2)$. $\qquad\qquad\qquad \square$

## G.3   Proof of Lemma 21

*proof.* Recall that Lemma 21 was stated and used in Appendix F with respect to Algorithm 3, therefore all events and variables in this subsection refer to those defined with respect to Algorithm 3. We will prove a more general result in Lemma 26.

**Lemma 26.** *Let $x_t, u_t$ respectively be the position and control of $C^{\mathrm{alg}}$ (the controller of Algorithm 3) at time $t$ starting at $x_0 = 0$. Define $G_i = (x_0, u_0, ..., x_{i-1}, u_{i-1})$. For constant $\gamma > 0$, define $S'_t$ as*

$$S'_t = \left\{ i < t : u_i = u_i^{\mathrm{safeU}} \text{ and } \mathbb{P}(u_i = u_i^{\mathrm{safeU}} \mid G_i, E) \geq \gamma \right\}, \tag{109}$$

*where $E$ is the event defined in Appendix F. Then under Assumptions 1–8 and for sufficiently large $T$, with probability $1 - o_T(1/T)$,*

$$\max_{s \in [0:s_e]} \epsilon_s \sqrt{|S'_{T_s}|} = \tilde{O}_T(1),$$

*where $\epsilon_s$ is from Algorithm 3.*

The proof of Lemma 26 can be found in Appendix G.4. We will now prove that Lemma 21 is a direct consequence of Lemma 26. By Equation (92), we have that for all $i$,

$$\mathbb{P}\left(u_i = u_i^{\mathrm{safeU}} \mid G_i, E\right) \geq \frac{\epsilon_{\mathrm{A9}}}{2}.$$

Therefore, we have that

$$\left\{ i < t : u_i = u_i^{\mathrm{safeU}} \text{ and } \mathbb{P}(u_i = u_i^{\mathrm{safeU}} \mid G_i, E) \geq \frac{\epsilon_{\mathrm{A9}}}{2} \right\} = \left\{ i < t : u_i = u_i^{\mathrm{safeU}} \right\}. \tag{110}$$

Lemma 26 for $\gamma = \frac{\epsilon_{\mathrm{A9}}}{2}$ gives that with probability $1 - o_T(1/T)$,

$$\max_{s \in [0:s_e]} \epsilon_s \cdot \sqrt{\left| \left\{ i < t : u_i = u_i^{\mathrm{safeU}} \text{ and } \mathbb{P}(u_i = u_i^{\mathrm{safeU}} \mid G_i, E) \geq \frac{\epsilon_{\mathrm{A9}}}{2} \right\} \right|} = \tilde{O}_T(1). \tag{111}$$

Combining Equation (110) and Equation (111) gives that with probability $1 - o_T(1/T)$,

$$\max_{s \in [0:s_e]} \epsilon_s \sqrt{\left| \left\{ i < t : u_i = u_i^{\mathrm{safeU}} \right\} \right|} = \tilde{O}_T(1),$$

which is the desired result of Lemma 21. $\qquad\square$

## G.4   Proof of Lemma 26

Lemma 26 is stated above to be used in Appendix F with respect to Algorithm 3, therefore all events and variables in this subsection refer to those defined for Algorithm 3.

*proof.* The first step of the proof will be to prove that conditional on event $E$ for all $i \geq T_0$,

$$u_i^{\mathrm{safeU}} = -\frac{a^*}{b^*} x_i + \frac{D_{\mathrm{U}} + e_i}{b^*}, \tag{112}$$

where $|e_i| = \tilde{O}_T(\nu_T)$. Let $s_i = \lfloor \log_2(i \nu_T^2) \rfloor$. Recall that $u_i^{\mathrm{safeU}}$ is the largest $u$ such that

$$\max_{\|\theta - \hat{\theta}_{s_i}\|_\infty \leq \epsilon_{s_i}} a x_i + b u \leq D_{\mathrm{U}}.$$

69

For sufficiently large $T$ and conditional on event $E$,

$$\epsilon_{s_i} = \tilde{O}_T(\nu_T) \leq \min(a^*, b^*) - \epsilon_{s_i} \leq \min(\hat{a}_{s_i}, \hat{b}_{s_i}).$$

This implies that $\hat{a}_{s_i} - \epsilon_{s_i} \geq 0$, giving the following equations of casework for $u_i^{\text{safeU}}$:

$$u_i^{\text{safeU}} = \begin{cases} \frac{D_U - (\hat{a}_{s_i} + \epsilon_{s_i})x_i}{\hat{b}_{s_i} - \epsilon_{s_i}}, & \text{if } x_i \geq 0 \text{ and } (\hat{a}_{s_i} + \epsilon_{s_i})x_i \geq D_U \\ \frac{D_U - (\hat{a}_{s_i} + \epsilon_{s_i})x_i}{\hat{b}_{s_i} + \epsilon_{s_i}}, & \text{if } x_i \geq 0 \text{ and } (\hat{a}_{s_i} + \epsilon_{s_i})x_i \leq D_U \\ \frac{D_U - (\hat{a}_{s_i} - \epsilon_{s_i})x_i}{\hat{b}_{s_i} + \epsilon_{s_i}}, & \text{if } x_i \leq 0 \end{cases} \tag{113}$$

which implies

$$u_i^{\text{safeU}} = \begin{cases} \frac{D_U - a^* x_i}{b^*} \frac{b^*}{\hat{b}_{s_i} - \epsilon_{s_i}} + \frac{a^* x_i - (\hat{a}_{s_i} + \epsilon_{s_i})x_i}{\hat{b}_{s_i} - \epsilon_{s_i}}, & \text{if } x_i \geq 0 \text{ and } (\hat{a}_{s_i} + \epsilon_{s_i})x_i \geq D_U \\ \frac{D_U - a^* x_i}{b^*} \frac{b^*}{\hat{b}_{s_i} + \epsilon_{s_i}} + \frac{a^* x_i - (\hat{a}_{s_i} + \epsilon_{s_i})x_i}{\hat{b}_{s_i} + \epsilon_{s_i}}, & \text{if } x_i \geq 0 \text{ and } (\hat{a}_{s_i} + \epsilon_{s_i})x_i \leq D_U \\ \frac{D_U - a^* x_i}{b^*} \frac{b^*}{\hat{b}_{s_i} + \epsilon_{s_i}} + \frac{a^* x_i - (\hat{a}_{s_i} - \epsilon_{s_i})x_i}{\hat{b}_{s_i} + \epsilon_{s_i}}, & \text{if } x_i \leq 0 \end{cases} \tag{114}$$

which implies

$$u_i^{\text{safeU}} = \begin{cases} \frac{D_U - a^* x_i}{b^*} + \frac{b^* - \hat{b}_{s_i} + \epsilon_{s_i}}{\hat{b}_{s_i} - \epsilon_{s_i}} \frac{D_U - a^* x_i}{b^*} + \frac{(a^* - \hat{a}_{s_i} - \epsilon_{s_i})x_i}{\hat{b}_{s_i} - \epsilon_{s_i}}, & \text{if } x_i \geq 0 \text{ and } (\hat{a}_{s_i} + \epsilon_{s_i})x_i \geq D_U \\ \frac{D_U - a^* x_i}{b^*} + \frac{b^* - \hat{b}_{s_i} - \epsilon_{s_i}}{\hat{b}_{s_i} + \epsilon_{s_i}} \frac{D_U - a^* x_i}{b^*} + \frac{(a^* - \hat{a}_{s_i} - \epsilon_{s_i})x_i}{\hat{b}_{s_i} + \epsilon_{s_i}}, & \text{if } x_i \geq 0 \text{ and } (\hat{a}_{s_i} + \epsilon_{s_i})x_i \leq D_U \\ \frac{D_U - a^* x_i}{b^*} + \frac{b^* - \hat{b}_{s_i} - \epsilon_{s_i}}{\hat{b}_{s_i} + \epsilon_{s_i}} \frac{D_U - a^* x_i}{b^*} + \frac{(a^* - \hat{a}_{s_i} + \epsilon_{s_i})x_i}{\hat{b}_{s_i} + \epsilon_{s_i}}. & \text{if } x_i \leq 0. \end{cases} \tag{115}$$

Under event $E$, $|a^* - \hat{a}_{s_i}| \leq \epsilon_{s_i}$, $|b^* - \hat{b}_{s_i}| \leq \epsilon_{s_i}$, and $|x_i| = \tilde{O}_T(1)$, therefore in all three cases we have that $u_i^{\text{safeU}} = -\frac{a^*}{b^*}x_i + \frac{D_U + e_i}{b^*}$, for some $e_i$ satisfying

$$|e_i| = \tilde{O}_T(\epsilon_{s_i}) = \tilde{O}_T(\nu_T). \tag{116}$$

We now define

$$S_t'' = \left\{ i < t : u_i = u_i^{\text{safeU}} \text{ and } \mathbb{P}(u_i = u_i^{\text{safeU}} \mid G_i, E) \geq \gamma \text{ and } \mathbb{P}(E \mid G_i) \geq \frac{1}{2} \right\}.$$

**Lemma 27.** *Using the same notation and assumptions as in the proof of Lemma 26, for any constant $c < 1$,*

$$\mathbb{P}\left( \forall i \in [0 : t - 1], \mathbb{P}(E \mid G_i) \geq c \right) = 1 - o_T(1/T).$$

*proof.* Consider any fixed $i \in [0 : t-1]$. We will show that $\mathbb{P}\left( \mathbb{P}(E \mid G_i) \geq c \right) = 1 - o_T(1/T^2)$. Suppose this is not true, i.e. suppose that $\mathbb{P}\left( \mathbb{P}(E \mid G_i) \geq c \right) = 1 - \Omega_T(1/T^2)$, or equivalently that $\mathbb{P}\left( \mathbb{P}(\neg E \mid G_i) \geq 1 - c \right) = \Omega_T(1/T^2)$. Note that by the law of total expectation,

$$\mathbb{P}\left( \neg E \mid \mathbb{P}(\neg E \mid G_i) \geq 1 - c \right) = \mathbb{E}\left[ \mathbb{P}\left( \neg E \mid G_i, \mathbb{P}(\neg E \mid G_i) \geq 1 - c \right) \mid \mathbb{P}(\neg E \mid G_i) \geq 1 - c \right]$$
$$\geq \mathbb{E}\left[ 1 - c \right]$$
$$= 1 - c.$$

This implies that

$$\mathbb{P}(\neg E) = \mathbb{P}\Big(\neg E \mid \mathbb{P}(\neg E \mid G_i) \geq 1 - c\Big)\mathbb{P}\Big(\mathbb{P}(\neg E \mid G_i) \geq 1 - c\Big) = (1 - c)\Omega_T(1/T^2).$$

This would then imply that $\mathbb{P}(E) = 1 - \mathbb{P}(\neg E) = 1 - \Omega_T(1/T^2)$, which is a contradiction with the fact that $\mathbb{P}(E) = 1 - o_T(1/T^2)$. Therefore, we must have that for all fixed $i$,

$$\mathbb{P}\Big(\mathbb{P}(E \mid G_i) \geq c\Big) = 1 - o_T(1/T^2).$$

Taking a union bound gives that

$$\mathbb{P}\Big(\forall i \in [0 : t - 1], \mathbb{P}(E \mid G_i) \geq c\Big) \geq 1 - \sum_{i=0}^{t-1} (1 - \mathbb{P}(\mathbb{P}(E \mid G_i) \geq c)) = 1 - o_T(1/T),$$

which is exactly what we want to show. □

If $\forall i \in [0 : t - 1]$, $\mathbb{P}(E \mid G_i) \geq 1/2$, then $|S_t'| = |S_t''|$. Using Lemma 27 with $c = 1/2$,

$$\mathbb{P}\left(|S_t'| = |S_t''|\right) \geq \mathbb{P}(\forall i \in [0 : t - 1], \mathbb{P}(E \mid G_i) \geq 1/2) = 1 - o_T(1/T). \tag{117}$$

Therefore, if we can show that with probability $1 - o_T(1/T)$,

$$\max_{s \in [0:s_e]} \epsilon_s \sqrt{|S_{T_s}''|} = \tilde{O}_T(1), \tag{118}$$

then a union bound combining Equation (118) with Equation (117) gives that with probability $1 - o_T(1/T)$,

$$\max_{s \in [0:s_e]} \epsilon_s \sqrt{|S_{T_s}'|} = \tilde{O}_T(1),$$

which is our desired result. Therefore, the rest of this proof will focus on proving Equation (118).

Fix any $s \in [0 : s_e]$. We will use Lemma 23 with $S = S_{T_s}''$. Under event $E$, we have by Lemma 4 the following three equations:

$$(V_{T_s}^{S_{T_s}''})_{11} = \lambda + \sum_{i=0}^{T_s-1} x_i^2 \mathbb{1}_{i \in S_{T_s}''} \leq \lambda + |S_{T_s}''| B_x^2 \tag{119}$$

$$(V_{T_s}^{S_{T_s}''})_{22} = \lambda + \sum_{i=0}^{T_s-1} u_i^2 \mathbb{1}_{i \in S_{T_s}''} \leq \lambda + |S_{T_s}''| B_x^2 \tag{120}$$

$$(V_{T_s}^{S_{T_s}''})_{12}^2 = \left(\sum_{i=0}^{T_s-1} u_i x_i \mathbb{1}_{i \in S_{T_s}''}\right)^2. \tag{121}$$

We can now lower bound $(V_{T_s}^{S''_{T_s}})_{22}(V_{T_s}^{S''_{T_s}})_{11} - (V_{T_s}^{S''_{T_s}})_{12}^2$ for sufficiently large $T$ conditional on event $E$.

$$(V_{T_s}^{S''_{T_s}})_{22}(V_{T_s}^{S''_{T_s}})_{11} - (V_{T_s}^{S''_{T_s}})_{12}^2$$

$$= \left(\lambda + \sum_{i=0}^{T_s-1} u_i^2 1_{i \in S''_{T_s}}\right)\left(\lambda + \sum_{i=0}^{T_s-1} x_i^2 1_{i \in S''_{T_s}}\right) - \left(\left(\sum_{i=0}^{T_s-1} u_i x_i 1_{i \in S''_{T_s}}\right)^2\right)$$

$$\geq \left(\sum_{i=0}^{T_s-1} u_i^2 1_{i \in S''_{T_s}}\right)\left(\sum_{i=0}^{T_s-1} x_i^2 1_{i \in S''_{T_s}}\right) - \left(\left(\sum_{i=0}^{T_s-1} u_i x_i 1_{i \in S''_{T_s}}\right)^2\right)$$

$$= \sum_{\substack{i<j}}^{T_s-1} (u_i x_j - u_j x_i)^2 1_{i,j \in S''_{T_s}}$$

$$= \sum_{\substack{i<j}}^{T_s-1} \left(\left(-\frac{a^*}{b^*}x_i + \frac{D_U + e_i}{b^*}\right)x_j - \left(-\frac{a^*}{b^*}x_j + \frac{D_U + e_j}{b^*}\right)x_i\right)^2 1_{i,j \in S''_{T_s}} \quad \text{Equation (112)}$$

$$= \frac{1}{(b^*)^2}\sum_{\substack{i<j}}^{T_s-1} (D_U x_j + e_i x_j - D_U x_i - e_j x_i)^2 1_{i,j \in S''_{T_s}}$$

$$= \frac{1}{(b^*)^2}\sum_{\substack{i<j}}^{T_s-1} (x_j(D_U + e_i) - (D_U + e_j)x_i)^2 1_{i,j \in S''_{T_s}}$$

$$= \frac{1}{(b^*)^2}\sum_{j=0}^{T_s-1} X_j 1_{j \in S''_{T_s}}. \tag{122}$$

Above we defined the random variable $X_j$ as

$$X_j = \sum_{i=0}^{j-1}((D_U + e_j)x_i 1_{i \in S''_{T_s}} - (D_U + e_i)x_j 1_{i \in S''_{T_s}})^2$$

$$\leq |S''_j| 4(D_U + 1)^2 B_x^2, \qquad \text{Equation (116)} \tag{123}$$

where the last inequality holds by Lemma 4 and because $e_j \leq 1$ under event $E$ for sufficiently large $T$ by Equation (116). We need one last lemma to help lower bound the conditional expectation of $X_j$.

**Lemma 28.** *Using the same notation and assumptions as in the proof of Lemma 26 (and recall that $B_P$ is the upper bound on the density of the noise random variables), if $\mathbb{P}(u_j = u_j^{\text{safeU}} \mid G_j, E) \geq \gamma$ and $\mathbb{P}(E \mid G_j) \geq 1/2$, then $\text{Var}\left(w_{j-1} \mid G_j, E, u_j = u_j^{\text{safeU}}\right) \geq \frac{\gamma^2}{64 B_P}$.*

The proof of Lemma 28 can be found in Appendix G.8. By definition, $j \in S''_{T_s}$ implies three events: $\{u_j = u_j^{\text{safeU}}\}$, $\{\mathbb{P}(u_j = u_j^{\text{safeU}} \mid G_j, E) \geq \gamma\}$, and $\{\mathbb{P}(E \mid G_j) \geq 1/2\}$. Note that the second and third events are deterministic functions of $G_j$. Therefore in the algebra below, the information in $\{j \in S''_{T_s}\}$ that tells us that $\mathbb{P}(u_j = u_j^{\text{safeU}} \mid G_j, E) \geq \gamma$ and $\mathbb{P}(E \mid G_j) \geq 1/2$ will be absorbed into the conditioning on $G_j$ in the first equality, i.e.,

starting in the second line below, the $G_j$ being conditioned on should be understood to be one for which $\mathbb{P}(u_j = u_j^{\text{safeU}} \mid G_j, E) \geq \gamma$ and $\mathbb{P}(E \mid G_j) \geq 1/2$. For sufficiently large $T$,

$$\mathbb{E}[X_j \mid G_j, E, j \in S''_{T_s}]$$

$$= \mathbb{E}[X_j \mid G_j, E, u_j = u_j^{\text{safeU}}]$$

$$= \mathbb{E}\left[\sum_{i=0}^{j-1}((D_U + e_j)x_i 1_{i \in S''_{T_s}} - (D_U + e_i)x_j 1_{i \in S''_{T_s}})^2 \mid G_j, E, u_j = u_j^{\text{safeU}}\right]$$

$$= \sum_{i=0}^{j-1}\mathbb{E}[((D_U + e_j)x_i 1_{i \in S''_{T_s}} - (D_U + e_i)(ax_{j-1} + bu_{j-1})1_{i \in S''_{T_s}} - (D_U + e_i)w_{j-1}1_{i \in S''_{T_s}})^2 \mid G_j, E, u_j = u_j^{\text{safeU}}]$$

$$\geq \sum_{i=0}^{j-1}\text{Var}\left((D_U + e_j)x_i 1_{i \in S''_{T_s}} - (D_U + e_i)(ax_{j-1} + bu_{j-1})1_{i \in S''_{T_s}} - (D_U + e_i)w_{j-1}1_{i \in S''_{T_s}} \mid G_j, E, u_j = u_j^{\text{safeU}}\right)$$

$$= \sum_{i=0}^{j-1}\text{Var}\left((D_U + e_i)w_{j-1} - e_j x_i \mid G_j, E, u_j = u_j^{\text{safeU}}\right)1_{i \in S''_{T_s}}$$

$$= \sum_{i=0}^{j-1}(D_U + e_i)^2 \text{Var}\left(w_{j-1} - \frac{e_j x_i}{D_U + e_i} \,\Big|\, G_j, E, u_j = u_j^{\text{safeU}}\right)1_{i \in S''_{T_s}}$$

$$\geq \sum_{i=0}^{j-1}(D_U + e_i)^2 \left(\text{Var}\left(w_{j-1} \mid G_j, E, u_j = u_j^{\text{safeU}}\right) - \left|2\text{Cov}\left(\frac{e_j x_i}{D_U + e_i}, w_{j-1} \,\Big|\, G_j, E, u_j = u_j^{\text{safeU}}\right)\right|\right)1_{i \in S''_{T_s}}$$

$$\geq \sum_{i=0}^{j-1}(D_U + e_i)^2 \left(\text{Var}\left(w_{j-1} \mid G_j, E, u_j = u_j^{\text{safeU}}\right) - \frac{1}{\log(T)}\right)1_{i \in S''_{T_s}} \qquad \text{Suff large } T \text{ (see below)}$$

$$\geq (D_U - \tilde{O}_T(\nu_T))^2 \cdot |S''_j| \cdot \left(\text{Var}\left(w_{j-1} \mid G_j, E, u_j = u_j^{\text{safeU}}\right) - \frac{1}{\log(T)}\right) \qquad \text{Equation (116)}$$

$$= (D_U - \tilde{O}_T(\nu_T))^2 \cdot |S''_j| \cdot \left(\frac{\gamma^2}{64 B_P} - \frac{1}{\log(T)}\right) \qquad \mathbb{P}(u_j = u_j^{\text{safeU}} \mid G_j, E) \geq \gamma, \text{ Lemma 28}$$

$$\geq \frac{D_U^2 |S''_j| \gamma^2}{128 B_P}. \qquad \text{Suff large } T \qquad (124)$$

Note that we are able to divide by $D_U + e_i$ for sufficiently large $T$ by Equation (116). The for-sufficiently-large-$T$ bound on the covariance comes from the fact that under event $E$, we have $|w_{j-1}| = \tilde{O}_T(1)$ and $\frac{e_j x_i}{D_U + e_i} = \tilde{O}_T(\nu_T)$, and therefore for sufficiently large $T$ the covariance has magnitude less than $\frac{1}{2\log(T)}$.

We can now apply Lemma 24 to $\{X_i\}_{i=0}^{T_s - 1}$ with $n = T_s$, $S_n = S''_{T_s}$, $p = \frac{D_U^2 \gamma^2}{1024 B_P (D_U + 1)^2 B_x^2}$, $F_i = G_i$, $E^* = E$, and $c = \frac{D_U^2 \gamma^2}{128 B_P}$ (where Equations (123) and (124) imply Equations (97) and (98)). Because $D_U \leq \log^2(T)$, $B_x = \log^3(T)$, and $\gamma$ is a constant, this choice of $p$ is less than 1 for sufficiently large $T$.

Applying Lemma 24 gives that for sufficiently large $T$, conditional on event $E$ with

conditional probability $1 - o_T(1/T^2)$,

$$(V_{T_s}^{S_{T_s}''})_{22}(V_{T_s}^{S_{T_s}''})_{11} - (V_{T_s}^{S_{T_s}''})_{12}^2$$

$$\geq \frac{1}{(b^*)^2} \sum_{j=0}^{T_s-1} X_j \qquad\qquad\qquad\qquad \text{Equation (122)}$$

$$\geq \frac{1}{(b^*)^2} \frac{D_U^2 \gamma^2}{512 B_P} \left( \max\left( \left\lfloor p|S_{T_s}''| - \sqrt{|S_{T_s}''| \log(T)} \right\rfloor, 1 \right) - 1 \right)$$

$$\times \left( \max\left( \left\lfloor p|S_{T_s}''| - \sqrt{|S_{T_s}''| \log(T)} \right\rfloor, 1 \right) \right) \qquad \text{Lemma 24} \qquad (125)$$

Define $E_s'$ as the event that Equation (125) holds (therefore $\mathbb{P}(E_s' \mid E) = 1 - o_T(1/T^2)$). If $|S_{T_s}''| \geq 4 \log^2(T)/p^2$, then $\frac{p|S_{T_s}''|}{2} \geq \log(T)\sqrt{|S_{T_s}''|}$, and therefore

$$p|S_{T_s}''| - \log(T)\sqrt{|S_{T_s}''|} \geq \frac{p|S_{T_s}''|}{2} \geq 1. \qquad (126)$$

Therefore, conditional on $E \cap E_s' \cap \{|S_{T_s}''| \geq 4 \log^2(T)/p^2\}$,

$$(V_{T_s}^{S_{T_s}''})_{22}(V_{T_s}^{S_{T_s}''})_{11} - (V_{T_s}^{S_{T_s}''})_{12}^2$$

$$\geq \frac{1}{(b^*)^2} \frac{D_U^2 \gamma^2}{512 B_P} \left( \max\left( \left\lfloor p|S_{T_s}''| - \sqrt{|S_{T_s}''| \log(T)} \right\rfloor, 1 \right) - 1 \right)$$

$$\times \left( \max\left( \left\lfloor p|S_{T_s}''| - \sqrt{|S_{T_s}''| \log(T)} \right\rfloor, 1 \right) \right) \qquad \text{Equation (125)}$$

$$\geq \frac{1}{(b^*)^2} \frac{D_U^2 \gamma^2}{512 B_P} \left( \lfloor p|S_{T_s}''|/2 \rfloor - 1 \right) \left( \lfloor p|S_{T_s}''|/2 \rfloor \right) \qquad \text{Equation (126)}$$

$$= \tilde{\Omega}_T \left( |S_{T_s}''|^2 \right). \qquad (127)$$

Because $E \subseteq E_1$, we have by Equation (108) that $B_{T_s} = \tilde{O}_T(1)$ conditional on event $E \cap N$. Therefore, by Lemma 23 and Equations (119), (120), (127), and (108), we have conditional on event $E \cap E_s' \cap N \cap \{|S_{T_s}''| \geq 4 \log^2(T)/p^2\}$ and for sufficiently large $T$,

$$\epsilon_s^2 \leq \frac{\lambda + |S_{T_s}''| B_x^2 \tilde{O}_T(1)}{\tilde{\Omega}(|S_{T_s}''|^2)} \leq \tilde{O}_T \left( \frac{1}{|S_{T_s}''|} \right). \qquad (128)$$

Taking $E' = \cap_{s \in [0:s_e]} E_s'$, Equation (128) implies that conditional on $E \cap E' \cap N \cap \{|S_{T_s}''| \geq 4 \log^2(T)/p^2\}$,

$$\epsilon_s^2 |S_{T_s}''| \leq \tilde{O}_T(1). \qquad (129)$$

Under event $E$, because $E_2 \subseteq E$, $\epsilon_s = \tilde{O}_T(\nu_T)$. Therefore, conditional on $E \cap E' \cap N \cap \{|S_{T_s}''| < 4 \log^2(T)/p^2\}$,

$$\epsilon_s^2 |S_{T_s}''| \leq \epsilon_s^2 \cdot 4 \log^2(T)/p^2 = \tilde{O}_T(\nu_T^2) = \tilde{O}_T(1). \qquad (130)$$

Because Equations (129) and (130) hold for all $s \in [0 : s_e]$, the right hand sides do not depend on $s$, and the equations hold almost surely, these two equations together imply that conditional on $E \cap E' \cap N$,

$$\max_{s \in [0:s_e]} \epsilon_s^2 |S_{T_s}''| = \tilde{O}_T(1).$$

74

Because $\mathbb{P}(E) \geq 1 - o_T(1/T)$, $\mathbb{P}(N) \geq 1 - o_T(1/T)$, and $\mathbb{P}(E' \mid E) \geq 1 - \sum_{s=0}^{s_e} \mathbb{P}(E'_s \mid E) = 1 - o_T(1/T)$, by a union bound we can conclude that with probability $1 - o_T(1/T)$,

$$\max_{s \in [0:s_e]} \epsilon_s^2 |S''_{T_s}| = \tilde{O}_T(1).$$

This completes the proof of Equation (118), and therefore completes the proof of this lemma.

$\square$

## G.5 Proof of Lemma 23

Recall that Lemma 23 applies to all algorithms as defined in the lemma statement, and therefore this lemma is not specific to a previous appendix section.

*proof.* First, we restate the theorem from Abbasi-Yadkori and Szepesvári [2011] in the notation and setup of this paper.

**Lemma 29** (Restatement of Theorem 1 in Abbasi-Yadkori and Szepesvári [2011])**.** *Let $\theta^* \in \mathbb{R}^2$ and $C$ be a controller. For $t \in [0 : T - 1]$, define $z_t = (x_t, C(x_t))$ and $x_{t+1} = \theta^* \cdot z_t + w_t$ where $w_t \sim_{i.i.d.} \mathcal{D}$ and $\mathcal{D}$ a subgaussian distribution with mean 0 and variance 1, and $\|\theta^*\|_2 \leq \bar{a}^2 + \bar{b}^2$. Define $V_t = \lambda I + \sum_{s=0}^{t-1} z_s z_s^\top$, $Z_t$ as the matrix where row $i \in [1 : t]$ is $z_{i-1}^\top$, and $X_t$ as the matrix where row $i \in [1 : t]$ is $x_i$. Finally, let $\hat{\theta}_t = (Z_t^\top Z_t + \lambda I)^{-1} Z_t^\top X_t$ and $\Delta_t = \hat{\theta} - \theta^*$. Then with probability $1 - o_T(1/T^2)$, for all $1 \leq t \leq T$.*

$$\text{Tr}(\Delta_t^\top V_t \Delta_t) \leq B_t^2, \tag{131}$$

*where $B_t = \alpha\sqrt{\log(\det(V_t)) + \log(\lambda^2) + 2\log(T^2)} + \sqrt{\lambda}(\bar{a}^2 + \bar{b}^2)$ and $\alpha$ satisfies $\mathbb{E}_{w \sim \mathcal{D}}[\exp(\gamma w)] \leq \exp(\gamma^2 \alpha^2/2)$ for any $\gamma \in \mathbb{R}$.*

Now define $V_t^{S^c} = \lambda I + \sum_{s=0}^{t-1} z_s z_s^\top \mathbb{1}_{s \notin S}$. Then by Lemma 29,

$$B_t \geq \text{Tr}(\Delta_t^\top V_t \Delta_t) = \text{Tr}(\Delta_t^\top (V_t^S + V_t^{S^c}) \Delta_t) = \text{Tr}(\Delta_t^\top V_t^S \Delta_t) + \text{Tr}(\Delta_t^\top V_t^{S^c} \Delta_t).$$

Because both traces are non-negative, this implies that $\text{Tr}(\Delta_t^\top V_t^S \Delta_t) \leq B_t^2$. Suppose $\Delta_t = (\Delta_{ta}, \Delta_{tb})$. Then expanding the trace gives that

$$(V_t^S)_{11}\Delta_{ta}^2 + (V_t^S)_{22}\Delta_{tb}^2 + 2\Delta_{ta}\Delta_{tb}(V_t^S)_{12} \leq B_t^2.$$

The left side of the above equation is a quadratic in $\Delta_{tb}$, with minimum occurring at $\Delta_{tb} = \frac{-\Delta_{ta}(V_t^S)_{12}}{(V_t^S)_{22}}$. Therefore, plugging this in gives the following inequality.

$$(V_t^S)_{11}\Delta_{ta}^2 - \frac{\Delta_{ta}^2(V_t^S)_{12}^2}{(V_t^S)_{22}} \leq B_t^2.$$

Simplifying, we have the desired result that $\Delta_{ta}^2 \leq \frac{(V_t^S)_{22}}{(V_t^S)_{11}(V_t^S)_{22} - (V_t^S)_{12}^2} B_t^2$. The proof follows symmetrically for $\Delta_{tb}$.

$\square$

## G.6 Proof of Lemma 24

*proof.* By the law of total expectation, for all $k \in [0, n-1]$,

$$c|S_k| \leq \mathbb{E}\left[X_k \mid F_k, E^*, k \in S_n\right] \quad \text{Eq (97)}$$

$$= \mathbb{E}\left[X_k \mid F_k, E^*, k \in S_n, X_k \leq \frac{c|S_k|}{2}\right] \mathbb{P}\left(X_k \leq \frac{c|S_k|}{2} \mid F_k, E^*, k \in S_n\right)$$

$$+ \mathbb{E}\left[X_k \mid X_k > \frac{c|S_k|}{2}, F_k, E^*, k \in S_n\right] \mathbb{P}\left(X_k > \frac{c|S_k|}{2} \mid F_k, E^*, k \in S_n\right)$$

$$\leq \frac{c|S_k|}{2} + \frac{c|S_k|}{2p}\mathbb{P}\left(X > \frac{c|S_k|}{2} \mid F_k, E^*, k \in S_n\right). \quad \text{Eq (98)}$$

For $i \in [0 : |S_n| - 1]$, define $\kappa_i$ as the $(i+1)$th smallest index in the set $S_n$. This implies that $|S_{\kappa_i}| = i$ and $\kappa_i \in S_n$. By Equation (98), for all $k$,

$$\mathbb{P}\left(X_k \geq \frac{c|S_k|}{2} \mid F_k, E^*, k \in S_n, \kappa_{|S_k|} = k\right) = \mathbb{P}\left(X_k \geq \frac{c|S_k|}{2} \mid F_k, E^*, k \in S_n\right) \geq \frac{c|S_k|/2}{c|S_k|/2p} = p. \quad (132)$$

Note that the first equality comes from the fact that by definition, $\kappa_{|S_k|} = k$ if $k \in S_n$, and $|S_k|$ is a deterministic function of $F_k$.

Let $A_0, A_1, ..., A_{n-1}$ be a sequence of i.i.d. Bernoulli random variables with probability $p$ of being 1 that are independent of all other random variables in this lemma, including $E^*, S_n, X_i, F_i$ for all $i$. For $i \in [0 : n-1]$, define the random variable $A'_i$ as

$$A'_i = \begin{cases} 1_{X_{\kappa_i} \geq \frac{c \cdot i}{2}} & \text{if } i \leq |S_n| - 1 \\ A_i & \text{otherwise.} \end{cases}$$

Define $F_i^A := F_{\kappa_{\min(i, |S_n|-1)}} \cup \{A_0, ..., A_{i-1}\}$. By Equation (132), we have that for all $i$,

$$\mathbb{P}\left(A'_i = 1 \mid F_i^A, E^*, i \leq |S_n| - 1\right)$$

$$= \sum_{k=0}^{n-1} \mathbb{P}\left(A'_i = 1 \mid F_i^A, E^*, i \leq |S_n| - 1, \kappa_i = k\right) \mathbb{P}\left(\kappa_i = k \mid F_i^A, E^*, i \leq |S_n| - 1\right) \quad \text{LoTE}$$

$$= \sum_{k=0}^{n-1} \mathbb{P}\left(X_{\kappa_i} \geq \frac{c \cdot i}{2} \mid F_i^A, E^*, i \leq |S_n| - 1, \kappa_i = k\right) \mathbb{P}\left(\kappa_i = k \mid F_i^A, E^*, i \leq |S_n| - 1\right)$$

$$= \sum_{k=0}^{n-1} \mathbb{P}\left(X_{\kappa_i} \geq \frac{c \cdot |S_{\kappa_i}|}{2} \mid F_i^A, E^*, i \leq |S_n| - 1, \kappa_i = k\right) \mathbb{P}\left(\kappa_i = k \mid F_i^A, E^*, i \leq |S_n| - 1\right)$$

$$= \sum_{k=0}^{n-1} \mathbb{P}\left(X_{\kappa_i} \geq \frac{c \cdot |S_{\kappa_i}|}{2} \mid F_{\kappa_i}, E^*, \kappa_i \in S_n, \kappa_i = k\right) \mathbb{P}\left(\kappa_i = k \mid F_i^A, E^*, i \leq |S_n| - 1\right)$$

$$= \sum_{k=0}^{n-1} \mathbb{P}\left(X_k \geq \frac{c \cdot |S_k|}{2} \mid F_k, E^*, k \in S_n, \kappa_i = k\right) \mathbb{P}\left(\kappa_i = k \mid F_i^A, E^*, i \leq |S_n| - 1\right)$$

$$= \sum_{k=0}^{n-1} \mathbb{P}\left(X_k \geq \frac{c \cdot |S_k|}{2} \mid F_k, E^*, k \in S_n, \kappa_{|S_k|} = k\right) \mathbb{P}\left(\kappa_i = k \mid F_i^A, E^*, i \leq |S_n| - 1\right) \quad i = |S_{\kappa_i}| = |S_k|$$

$$\geq \sum_{k=0}^{n-1} p \cdot \mathbb{P}\left(\kappa_i = k \mid F_i^A, E^*, i \leq |S_n| - 1\right) \quad \text{Eq (132)}$$

$$= p, \quad (133)$$

76

and

$$\mathbb{P}\left(A'_i = 1 \mid F_i^A, E^*, i > |S_n| - 1\right)$$
$$= \mathbb{P}\left(A_i = 1 \mid F_i^A, E^*, i > |S_n| - 1\right) \qquad \text{Independence of } A_i$$
$$= p. \tag{134}$$

Putting together Equations (133) and (134) and the Law of Total Probability,

$$\mathbb{P}\left(A'_i = 1 \mid F_i^A, E^*\right)$$
$$= \mathbb{P}\left(A'_i = 1 \mid F_i^A, E^*, i \le |S_n| - 1\right) \mathbb{P}\left(i \le |S_n| - 1 \mid F_i^A, E^*\right)$$
$$\quad + \mathbb{P}\left(A'_i = 1 \mid F_i^A, E^*, i > |S_n| - 1\right) \mathbb{P}\left(i > |S_n| - 1 \mid F_i^A, E^*\right)$$
$$\ge p. \qquad \text{Eqs (133) and (134).} \tag{135}$$

Because $A'_i$ is a deterministic function of $F_{i+1}^A$ and $F_i^A \subseteq F_{i+1}^A$, Equation (135) implies that $M_k = \sum_{i=0}^{k-1}(A'_i - p)$ is a submartingale conditional on $E^*$ with increments bounded in magnitude by 1. For any non-random $m \in [1:n]$, the Azuma–Hoeffding Inequality therefore gives that

$$\mathbb{P}\left(\sum_{i=0}^{m-1}(A'_i - p) \ge -\log(T)\sqrt{m} \;\middle|\; E^*\right) \ge 1 - e^{-\log^2(T)m/(2m)} = 1 - o_T(1/T^3).$$

Taking a union bound over all $m \in [1:n]$ (because $n \le T$), we have that

$$\mathbb{P}\left(\forall m \in [1:n], \sum_{i=0}^{m-1} A'_i \ge pm - \log(T)\sqrt{m} \;\middle|\; E^*\right) \ge 1 - o_T(1/T^2).$$

Define $E'$ as the event that for all $m \in [1:n]$, $\sum_{i=0}^{m-1} A'_i \ge pm - \log(T)\sqrt{m}$. Because $|S_n| \in [0, n]$, we must have that conditional on event $E'$,

$$\sum_{i=0}^{|S_n|-1} A'_i \ge p|S_n| - \log(T)\sqrt{|S_n|}. \tag{136}$$

Therefore, conditional on event $E'$, we have

$$\sum_{j=0}^{n-1} X_j \ge \sum_{j=0, j \in S_n}^{n-1} X_j \qquad\qquad X_j \ge 0$$

$$\ge \sum_{i=0}^{|S_n|-1} \frac{c \cdot i}{2} \cdot A'_i \qquad\qquad \text{Def of } A'_i$$

$$\ge \frac{c}{2} \sum_{k=0}^{\max(\lfloor p|S_n| - \log(T)\sqrt{|S_n|}\rfloor, 1) - 1} k. \qquad\qquad \text{Eq (136)}$$

$$= \frac{c}{4} \left(\max(\lfloor p|S_n| - \log(T)\sqrt{|S_n|}\rfloor, 1)\right) \left(\max(\lfloor p|S_n| - \log(T)\sqrt{|S_n|}\rfloor, 1) - 1\right)$$

Because we already showed that $\mathbb{P}(E' \mid E^*) \ge 1 - o_T(1/T^2)$, this is the desired result. $\qquad\square$

## G.7 Proof of Lemma 25

Recall that Lemma 25 was stated to be used in Appendix C with respect to Algorithm 2, therefore all events and variables in this subsection refer to those defined with respect to Algorithm 2.

*proof.* Define $A_i = 1_{|x_i| \leq \frac{1}{\log(T)}}$. Recall that $x_i = a^* x_{i-1} + b^* u_{i-1} + w_{i-1}$, where $x_{i-1}$ and $u_{i-1}$ are respectively the position and control at time $t = i - 1$. The probability that $A_i$ is equal to 1 is the probability that $w_{i-1} \in [-(a^* x_{i-1} + b^* u_{i-1}) - \frac{1}{\log(T)}, -(a^* x_{i-1} + b^* u_{i-1}) + \frac{1}{\log(T)}]$. Because $\mathcal{D}$ has a bounded density function (bounded by $B_P$) as assumed in Assumption 3, the conditional probability given $G_i$ is at most $\frac{2B_P}{\log(T)}$. Therefore, we have that

$$\mathbb{P}(A_i = 1 \mid G_i) \leq \frac{2B_P}{\log(T)}.$$

Therefore, $M_j = \sum_{i=0}^{j-1} (A_i - \frac{2B_P}{\log(T)})$ is a submartingale with differences bounded in magnitude by $\max(1, \frac{2B_P}{\log(T)}) \leq 1$ for sufficiently large $T$. By Azuma–Hoeffding's inequality, with probability $1 - o_T(1/T^3)$,

$$M_j \leq \log(T)\sqrt{j}.$$

Define $E_{\text{L25}}^j$ as the event that this bound on $M_j$ holds. By construction of $M_j$, under event $E_{\text{L25}}^j$,

$$\left| \left\{ i < j : |x_i| \leq \frac{1}{\log(T)} \right\} \right| = \sum_{i=0}^{j-1} A_i \leq \frac{2jB_P}{\log(T)} + \log(T)\sqrt{j} \leq \frac{4jB_P}{\log(T)}$$

for $j \geq \log^8(T)$ assuming $T$ is large enough that $\log^2(T) \geq \frac{1}{2B_P}$. As long as $\log(T) \geq 8B_P$, this implies that under event $E_{\text{L25}}^j$,

$$\left| \left\{ i < j : |x_i|^2 \geq \frac{1}{\log^2(T)} \right\} \right| \geq j - \frac{4jB_P}{\log(T)} \geq \frac{j}{2}.$$

Finally, we can conclude that under event $E_{\text{L25}}^j$,

$$\sum_{i=0}^{j-1} x_i^2 \geq \frac{j}{2\log^2(T)}.$$

We have shown that Equation (104) holds for any fixed $j$ under event $E_{\text{L25}}^j$ for sufficiently large $T$. Therefore, the same result holds for all $j \geq \log^8(T)$ under event $E_{\text{L25}} = \cap_{j \geq \log^8(T)} E_{\text{L25}}^j$. By a union bound and because $\mathbb{P}(E_{\text{L25}}^j) = 1 - o_T(1/T^3)$ for all $j$, we have that $\mathbb{P}(E_{\text{L25}}) = 1 - o_T(1/T^2)$. $\qquad\square$

## G.8 Proof of Lemma 28

Recall that Lemma 28 is defined to be used in Appendix F with respect to Algorithm 3, therefore all events and variables in this subsection refer to those defined with respect to Algorithm 3.

*proof.* By assumption of this lemma,

$$\mathbb{P}(u_j = u_j^{\text{safeU}}, E \mid G_j) = \mathbb{P}(u_j = u_j^{\text{safeU}} \mid G_j, E)\mathbb{P}(E \mid G_j) \geq \frac{\gamma}{2}. \tag{137}$$

We also note the following result:

**Lemma 30.** *For any event $E^*$ such that $\mathbb{P}(E^*) > 0$,*

$$\operatorname*{Var}_{w \sim \mathcal{D}}(w \mid E^*) \geq \frac{\mathbb{P}(E^*)^2}{16B_P^2}$$

*proof.* First, we will show that any continuous distribution $\mathcal{D}'$ with density function bounded by $B$ must have variance at least $\frac{1}{16B^2}$. Let $f_{\mathcal{D}'}$ be the probability density function of $\mathcal{D}'$. First, we can assume WLOG that $\mathcal{D}'$ has mean 0 (this is without loss of generality because variance is invariant to shifts in mean). If $\mathcal{D}'$ has mean 0, then by the law of total expectation

$$\mathbb{E}_{x \sim \mathcal{D}'}[x \mid x \geq 0]\mathbb{P}_{x \sim \mathcal{D}'}(x \geq 0) = -\mathbb{E}_{x \sim \mathcal{D}'}[x \mid x \leq 0]\mathbb{P}_{x \sim \mathcal{D}'}(x \leq 0).$$

Note that we can have non-strict inequalities because $\mathcal{D}'$ is continuous. Furthermore, either $\mathbb{P}_{x \sim \mathcal{D}'}(x \leq 0) \geq 1/2$ or $\mathbb{P}_{x \sim \mathcal{D}'}(x \geq 0) \geq 1/2$. Because variance is invariant to multiplying by $-1$, we can assume WLOG that $\mathbb{P}_{x \sim \mathcal{D}'}(x \geq 0) \geq 1/2$. If $\mathbb{P}_{x \sim \mathcal{D}'}(x \geq 0) \geq 1/2$ then $\int_0^\infty f_{\mathcal{D}'}(x)dx \geq 1/2$. Define $f^*(x) = \frac{1}{2B}$ for $x \in [0, B]$ and $f^*(x) = 0$ otherwise. Note that $f = f^*$ achieves the minimum possible value of $\int_0^\infty x \cdot f(x)dx$ subject to the constraints $\int_0^\infty f(x)dx \geq 1/2$ and $0 \leq f(x) \leq B$ for all $x$. This is because $f^*$ puts as much weight as possible close to 0 without violating the bounded by $B$ constraint. Furthermore, any $f$ such that $\int_{1/2B}^\infty f(x)dx > 0$ puts non-0 weight on values of $x$ greater than $B$ and therefore has a larger value of $\int_0^\infty x \cdot f(x)dx$ than $f^*$. Using this, we have that

$$\mathbb{E}_{x \sim \mathcal{D}'}[x \mid x \geq 0]\mathbb{P}_{x \sim \mathcal{D}'}(x \geq 0) = \int_0^\infty x \cdot f_{\mathcal{D}'}(x)dx \geq \int_0^{1/2B} x \cdot Bdx = \frac{1}{8B}.$$

Therefore, we must have (again by the law of total expectation) that

$$\mathbb{E}_{x \sim \mathcal{D}'}[|x|] = \mathbb{E}_{x \sim \mathcal{D}'}[x \mid x \geq 0]\mathbb{P}_{x \sim \mathcal{D}'}(x \geq 0) - \mathbb{E}_{x \sim \mathcal{D}'}[x \mid x \leq 0]\mathbb{P}_{x \sim \mathcal{D}'}(x \leq 0) \geq \frac{1}{4B}.$$

By Jensen's inequality,

$$\operatorname{Var}_{x \sim \mathcal{D}'}(x) = \mathbb{E}_{x \sim \mathcal{D}'}[x^2] = \mathbb{E}_{x \sim \mathcal{D}'}[|x|^2] \geq \mathbb{E}_{x \sim \mathcal{D}'}[|x|]^2 \geq \frac{1}{16B^2}.$$

We have therefore shown that any continuous distribution $\mathcal{D}'$ with probability density function $f$ such that $f(x) \leq B$ for all $x$ must have variance at least $\frac{1}{16B^2}$.

We know that the conditional distribution of $w$ given $E^*$ has a probability density function that is bounded by $\frac{B_P}{\mathbb{P}(E^*)}$. Therefore, we must have that $\operatorname{Var}(w \mid E^*) \geq \frac{\mathbb{P}(E^*)^2}{16B_P^2}$. $\qquad\square$

Recall that $w_{j-1}$ is independent of $G_j$. Therefore, $\text{Var}\left(w_{j-1} \mid G_j, E, u_j = u_j^{\text{safeU}}\right)$ is simply the variance of $w_{j-1}$ conditional on an event that has probability $\mathbb{P}(E, u_j = u_j^{\text{safeU}} \mid G_j)$. Therefore, we can apply Lemma 30 and Equation (137) to get that for some event $E'$ such that $\mathbb{P}(E') \geq \gamma/2$,

$$\text{Var}\left(w_{j-1} \mid G_j, E, u_j = u_j^{\text{safeU}}\right) = \text{Var}\left(w_{j-1} \mid E'\right) \geq \frac{\gamma^2}{64 B_P^2}.$$

$\square$

# H  Feasibility and Boundary Proofs

## H.1  Relaxation of Assumption 1

The assumption that $\underline{a}, \underline{b} > 0$ in Assumption 1 can actually be dropped under Assumptions 2 and 3. Informally, this is because the controller $C^{\text{init}}$ can be used for $\log^{10}(T)$ steps to, with high probability, obtain an estimate $\hat{\theta}$ such that $\|\hat{\theta} - \theta^*\|_\infty \leq \frac{1}{\log(T)}$ (by the same logic as in Lemma 2). Therefore, we could include an initial phase in every algorithm that does $\log^{10}(T)$ steps of initial exploration and then replaces $\Theta$ with $\Theta' = \{\theta : \|\theta - \hat{\theta}\|_\infty \leq \frac{1}{\log(T)}\}$, and this $\Theta'$ will satisfy $\underline{a}', \underline{b}' > 0$ for sufficiently large $T$ because $a^* > 0$. However, to simplify the algorithms and proofs we will assume that the initial uncertainty set $\Theta$ is small enough that this is unnecessary. Note that this assumption of sufficiently small bounded initial uncertainty appears in other safe LQR literature such as Li et al. [2021].

## H.2  Discussion on Assumption 2

To better understand Assumption 2, consider the case of bounded noise and constant boundaries as in Li et al. [2021], Dean et al. [2019]. In this case, to satisfy Assumption 2, it is sufficient to replace the $\forall x \in \left[D_{\text{L}}^{\mathbb{E}[x]} + F_{\mathcal{D}}^{-1}(\frac{1}{T^4}), D_{\text{U}}^{\mathbb{E}[x]} + F_{\mathcal{D}}^{-1}(1 - \frac{1}{T^4})\right]$ with $\forall x \in [D_{\text{L}}^{\mathbb{E}[x]} - \bar{w}, D_{\text{U}}^{\mathbb{E}[x]} + \bar{w}]$. Li et al. [2021] makes a similar assumption that there is an initial linear controller that satisfies this property. For the bounded noise case, Assumption 2 can be shown to be equivalent to an assumption on the size of the initial uncertainty set. Let $C^{\text{init}}(x_t) = -\frac{a}{b} x_t$ for some arbitrary $\theta \in \Theta$. When using this controller, the position and control at time $t$ (denoted $x_t$ and $u_t$ respectively) satisfy

$$|a^* x_t + b^* u_t| \leq |x_t| \left|a^* - \frac{ab^*}{b}\right| \leq |x_t| \left|a^* - a - \frac{(b^* - b)a}{b}\right| \leq \left(1 + \frac{a}{b}\right)|x_t|\text{size}(\Theta) \leq \left(1 + \frac{\bar{a}}{\underline{b}}\right)|x_t|\text{size}(\Theta).$$

This controller $C^{\text{init}}$ satisfies Assumption 2 under bounded noise if

$$\text{size}(\Theta) \leq \frac{\min(D_{\text{U}}^{\mathbb{E}[x]}, |D_{\text{L}}^{\mathbb{E}[x]}|) - \frac{\bar{b}}{\log(T)}}{\left|1 + \frac{\bar{a}}{\underline{b}}\right|(\|D^{\mathbb{E}[x]}\|_\infty + \bar{w})}.$$

Therefore, instead of assuming Assumption 2, it is sufficient to assume that $\text{size}(\Theta) \leq \frac{\min(D_{\text{U}}^{\mathbb{E}[x]}, |D_{\text{L}}^{\mathbb{E}[x]}|) - \frac{\bar{b}}{\log(T)}}{\left|1 + \frac{\bar{a}}{\underline{b}}\right|(\|D^{\mathbb{E}[x]}\|_\infty + \bar{w})}$, as the controller $C^{\text{init}}(x_t) = -\frac{a}{b} x_t$ satisfies Assumption 2. The bound

on size($\Theta$) does still depend on the end points of $\Theta$. As a sanity check, suppose $\|D^{\mathbb{E}[x]}\|_\infty = O_T(1)$ and $\bar{a}, \bar{b}, \frac{1}{\underline{b}} \leq c$ for some constant $c$. Then there exists a constant such that if size($\Theta$) is less than that constant, then Assumption 2 is satisfied for sufficiently large $T$.

## H.3   Assumptions Relationship to Infeasibility

In this section we briefly relate the assumptions we make to a notion of infeasibility. We begin with two formal definitions. The first is a formal definition of feasibility for our problem. The second is a property of a controller that is slightly stronger than regular safety.

**Definition 3** (Feasibility). *An initial uncertainty set of system dynamics $\Theta$ is* feasible *for boundary $D^{\mathbb{E}[x]}$ and trajectory length $T$ with probability $1 - \delta$ if there exists a controller $C$ that satisfies the following. For any $\theta^* \in \Theta$, if the true dynamics are $\theta^*$, then*

$$\mathbb{P}\left(\forall t < T : D_{\mathrm{L}}^{\mathbb{E}[x]} \leq a^* x_t + b^* C(H_t) \leq D_{\mathrm{U}}^{\mathbb{E}[x]}\right) \geq 1 - \delta.$$

**Definition 4** (Robust safety). *A controller $C$ is* robustly safe *for $T_0$ time steps for dynamics $\theta^*$ if the following holds for some known distribution $\rho$ with mean $0$ and constant variance $\eta^2 > 0$. If $s_t \overset{i.i.d.}{\sim} \rho$ and $u_t = C(H_t) + \frac{s_t}{\log(T)}$, then*

$$\mathbb{P}\left(\forall t \in [0, T_0 - 1] : D_{\mathrm{L}}^{\mathbb{E}[x]} \leq a^* x_t + b^* u_t \leq D_{\mathrm{U}}^{\mathbb{E}[x]}\right) \geq 1 - o_T(1/T^4).$$

**Proposition 11.** *The result of Theorem 1 hold without Assumption 2 if we assume access to a controller $C^{\mathrm{rs}}$ that is robustly safe for $\sqrt{T}$ steps. Similarly, the result of Theorem 2 holds without Assumption 2 if we assume access to a controller $C^{\mathrm{rs}}$ that is robustly safe for $T^{2/3}$ steps.*

*proof.* In the exploration phase of any of the three algorithms, instead of sampling $\phi_t$ from Rademacher distribution we can instead sample i.i.d. from $\rho$ and keep the rest of the algorithm the same. Then the robust feasibility implies that with probability $1 - o_T(1/T^4)$ the algorithm will be safe for the warm-up period of the first $\frac{1}{\nu_0^2}$ steps. We can then proof a variation of Lemma 2 that holds using the distribution $\rho$ instead of the Rademacher distribution. $\square$

By Definition 3, as $T$ approaches infinity, the existence of a robustly safe controller $C^{\mathrm{rs}}$ becomes intuitively equivalent to $\Theta$ being feasible for boundary $D^{\mathbb{E}[x]}$ with probability $1 - o_T(1/T^4)$. Therefore by Proposition 11, Assumption 2 is intuitively asymptotically equivalent to the assumption that the problem is feasible for dynamics $\Theta$ and that a controller that achieves feasibility is known.

# I   Generalizations

## I.1   Control Constraints

Our results focus on positional constraints, but we believe that our results with the same rates of regret will also hold with both positional and control constraints under some additional

assumptions. While we leave the formal derivations of results for control constraints to future work, we provide a brief discussion of how the algorithm and proofs from this paper could be extended to include control constraints.

First, we briefly mention how control constraints change the definitions and notation used. Control constraints would be of the form $D_{\mathrm{L}}^u \leq u_t \leq D_{\mathrm{U}}^u$ for all $t < T$ (for the rest of this section, we will refer to the expected-position constraints as $D^{\mathbb{E}[x]}$). We also define the function $K_{\mathrm{opt}}(\theta, T, D^{\mathbb{E}[x]}, D^u)$ as choosing the optimal parameter $K$ for a controller satisfying both the position constraints $D^{\mathbb{E}[x]}$ and the control constraints $D^u$. We also need the additional assumption that there exists a (non-empty) set of baseline controllers that can satisfy both the position and control constraints. Finally, we need to assume that the controller $C^{\mathrm{init}}$ satisfies both position and control constraints (i.e. an analogue of Assumption 2).

### I.1.1 Theorem 2 and Algorithm 2

We start with considering how Algorithm 2 would need to be modified with the addition of control constraints. The key idea behind Algorithm 2 satisfying the position constraints is that the algorithm sometimes uses controls $u_t^{\mathrm{safeU}}$ and $u_t^{\mathrm{safeL}}$ to enforce positional safety. However, in the presence of control constraints, we can no longer use the controls $u_t^{\mathrm{safeU}}$ and $u_t^{\mathrm{safeL}}$, as these controls may not satisfy the control constraints. The key modification of Algorithm 2 is to choose the controller $C_s^{\mathrm{alg}}$ in such a way that $C_s^{\mathrm{alg}}$ will satisfy a tighter positional constraint with respect to $D^{\mathbb{E}[x]'} = (D_L^{\mathbb{E}[x]} + \tilde{\Theta}_T(\epsilon_s), D_U^{\mathbb{E}[x]} - \tilde{\Theta}_T(\epsilon_s))$ for dynamics $\hat{\theta}_s$ and a tighter control constraint $D^{u'} = (D_L^u + \tilde{\Theta}_T(\epsilon_s), D_U^u - \tilde{\Theta}_T(\epsilon_s))$. In other words, choosing $C_s^{\mathrm{alg}} = C_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s, D^{\mathbb{E}[x]'}, D^{u'})}^{\hat{\theta}_s}$. Within each iteration of the safe exploitation phase, the algorithm then can directly use $C_s^{\mathrm{alg}}$. Because $\|\hat{\theta}_s - \theta^*\|_\infty \leq \tilde{O}_T(\epsilon_s)$ with high probability and this $C_s^{\mathrm{alg}}$ is chosen to satisfy the tighter position constraints $D^{\mathbb{E}[x]'}$ for dynamics $\hat{\theta}_s$, the controller $C_s^{\mathrm{alg}}$ will satisfy the true position constraints $D^{\mathbb{E}[x]}$ for dynamics $\theta^*$ with high probability. Because $C_s^{\mathrm{alg}}$ satisfies the tighter control constraints $D^{u'}$, with the additional assumption that the controller class is continuous, the controls used by $C_s^{\mathrm{alg}}$ under dynamics $\theta^*$ will also satisfy the control constraints with high probability.

Now we will briefly describe what additional results need to be proven in order for the modified version of Algorithm 2 described above to achieve the same regret rate of $\tilde{O}_T(T^{2/3})$ in the presence of control constraints. We will do this by analyzing each of the terms of regret from the proof of Theorem 2.

The regret term $R_0$, which is the regret from the warm-up period of the first $1/\nu_T^2$ steps, would have the same definition and the same regret bound of $\tilde{O}(T^{2/3})$ as in the analysis of Algorithm 2.

To bound the regret term $R_1$, we would need to show that $C_{\mathrm{alg}}^s$ as described above does not have much more expected cost than the true best controller, $C_{K_{\mathrm{opt}}(\theta^*, T, D^{\mathbb{E}[x]}, D^u)}^{\theta^*}$. This can be incorporated into an analogue of Assumption 7: assuming that for $\|\theta - \theta^*\|_\infty, \|D^u -$

$D^{u'}\|_\infty, \|D^{\mathbb{E}[x]} - D^{\mathbb{E}[x]'}\|_\infty$ all sufficiently small,

$$|\bar{J}(\theta^*, C^\theta_{K_{\mathrm{opt}}(\theta, T, D^{\mathbb{E}[x]'}, D^{u'})}, t) - \bar{J}(\theta^*, C^{\theta^*}_{K_{\mathrm{opt}}(\theta^*, T, D^{\mathbb{E}[x]}, D^u)}, t)|$$

$$= \tilde{O}_T\left(\|\theta - \theta^*\|_\infty + \|D^{\mathbb{E}[x]} - D^{\mathbb{E}[x]'}\|_\infty + \|D^u - D^{u'}\|_\infty + \frac{1}{T^2}\right).$$

This can be made into a new assumption on the baseline class of controllers that replaces Assumption 7.

We expect that the regret source $R_2$ (converting from expected regret to realized regret) will still be $\tilde{O}_T(\sqrt{T})$, as this was a result of a concentration inequality that will still apply.

Regret $R_3$ no longer exists as we no longer use the controls $u_t^{\mathrm{safeU}}$ or $u_t^{\mathrm{safeL}}$, and instead this source of regret is being incorporated into the chosen $C_s^{\mathrm{alg}}$ in regret term $R_1$.

To summarize, the main modification to the algorithm would be the choice of controller $C_s^{\mathrm{alg}}$, and the main change to the proof is moving the burden of bounding the regret term $R_3$ to the version of Assumption 7 described above that accounts for the tightened constraint arguments to $K_{\mathrm{opt}}$.

### I.1.2   Theorem 1 and Algorithm 3

In order to show a version of Theorem 1 that works for control constraints, Algorithm 3 would need the same modifications as described for Algorithm 2. Specifically, instead of using controls $u_t^{\mathrm{safeU}}$ and $u_t^{\mathrm{safeL}}$, the controller $C_s^{\mathrm{alg}}$ is chosen as $C_s^{\mathrm{alg}} = C^{\hat{\theta}_s}_{K_{\mathrm{opt}}(\hat{\theta}_s, T_s, D^{\mathbb{E}[x]'}, D^{u'})}$.

The main way that the proof of regret for Algorithm 3 differs from the regret for Algorithm 2 is that the proof for Algorithm 3 relies on the faster rate of convergence for $\hat{\theta}_s$ given by Lemma 21. Proving a form of Lemma 21 for the modified algorithm would be the main additional step in proving that $\tilde{O}_T(\sqrt{T})$ regret is possible with control constraints. As discussed in the proof sketch of Theorem 1, the proof of Lemma 21 comes from the fact that a constant fraction of the time, $u_t^{\mathrm{safeU}}$ is non-linear by an amount larger than a positive constant. The non-linearity of $u_t^{\mathrm{safeU}}$ occurs because enforcing safety constraint satisfaction requires non-linear controls. While the modified controller $C_s^{\mathrm{alg}}$ described in the previuos paragraph does not use the non-linear controls $u_t^{\mathrm{safeU}}$, $C_s^{\mathrm{alg}}$ must still be frequently non-linear in order to satisfy the safety constraints. Therefore, we expect that for noise distributions with large enough support, the modified Algorithm 3 will a constant fraction of the time use a control that is non-linear by a constant amount, which will give that $\epsilon_s$ decreases at a rate of $1/\sqrt{t}$.

## I.2   Higher Dimensions

This work focuses on the one-dimensional LQR setting, but many LQR applications have higher dimensional positions and controls. We leave the formal extension of our results to higher dimensions for future work, but discuss here when and how we believe our results will extend to higher dimensions. Suppose $x_t \in \mathbb{R}^n$ and $u_t \in \mathbb{R}^m$, which implies that the dynamics are a pair of matrices $\theta^* = (A^*, B^*)$ where $A^* \in \mathbb{R}^{n \times n}$ and $B^* \in \mathbb{R}^{n \times m}$. A natural extension of our constraints to higher dimensions is to consider a (origin-containing) polytopal constraint, i.e., the intersection of a finite number of half-spaces that contain the

origin. Specifically, we could consider constraints of the form $\Delta(A^* x_t + B^* u_t) \leq d$ where $\Delta \in \mathbb{R}^{k \times n}$ and $d \in \mathbb{R}^k$. This still has the interpretation as the expected position at each time is within the convex region $\{x \in \mathbb{R}^n : \Delta x \leq d\}$. Analogous to in Appendix I.1, we define the function $K_{\text{opt}}(\theta, T, \Delta, d)$ as choosing the optimal parameter $K$ for a controller satisfying the constraints $\Delta(A x_t + B u_t) \leq d$. Before talking about specific algorithms, we first note that we expect that the results of Lemmas 23 and 2 generalize directly to higher dimensions. This is necessary for all of our algorithmic results. Note that because the dynamics are matrices, the dynamics estimates will also be matrices denoted $\hat{\theta}_s$.

### I.2.1 Theorem 2 and Algorithm 2

In higher dimensions, Assumption 5 becomes slightly more complicated. Specifically, we define the truncated version of a controller $C$ in higher dimensions as using either control $C(x)$ if $C(x)$ would result in an expected position inside the convex safe region, and otherwise using the smallest magnitude control that takes the position in expectation to inside of the convex safe region. The other assumptions have direct higher dimensional counterparts.

The key modification of Algorithm 2 is to choose the controller $C_s^{\text{alg}}$ in such a way that $\Delta(\hat{A}_s x_t + \hat{B}_s C_s^{\text{alg}}(x_t)) \leq d - \tilde{\Theta}_T(\epsilon_s)$. In other words, choosing $C_s^{\text{alg}} = C_{K_{\text{opt}}(\hat{\theta}_s, T_s, \Delta, d - \tilde{\Theta}_T(\epsilon_s))}^{\hat{\theta}_s}$. Within each iteration of the main loop of Algorithm 2, the algorithm can directly use $C_s^{\text{alg}}$ without the need for $u_t^{\text{safeU}}$ or $u_t^{\text{safeL}}$. By this construction, $\Delta(\hat{A}_s x_t + \hat{B}_s C_s^{\text{alg}}(x_t)) \leq d - \tilde{O}_T(\epsilon)$. Because with high probability $\|\hat{\theta}_s - \theta^*\|_\infty \leq \tilde{O}_T(\epsilon_s)$, this will imply that $\Delta(A^* x_t + B^* C_s^{\text{alg}}(x_t)) \leq d$ with high probability. This in turn means that the algorithm will satisfy the constraints with high probability.

Analyzing the regret of this algorithm, the regret terms $R_0$, $R_1$, and $R_2$ stay the same as in the proof of Theorem 2. The regret term $R_3$ is no longer needed, as we no longer use controls $u_t^{\text{safeU}}$ or $u_t^{\text{safeL}}$. To bound the regret term $R_1$, we want to show that the cost of $C_{\text{alg}}^s$ is close to the cost of $C_{K_{\text{opt}}(\theta^*, T, \Delta, d)}^{\theta^*}$. Like we did in Appendix I.1, we need an analogue of Assumption 7, which is that for $\|\theta - \theta^*\|_\infty$ and $\|d - d'\|_\infty$ both sufficiently small,

$$
\left| \bar{J}(\theta^*, C_{K_{\text{opt}}(\theta, T, \Delta, d')}^{\theta}, t) - \bar{J}(\theta^*, C_{K_{\text{opt}}(\theta^*, T, \Delta, d)}^{\theta^*}, t) \right|
$$
$$
= \tilde{O}_T \left( \|\theta - \theta^*\|_\infty + \|d - d'\|_\infty + \frac{1}{T^2} \right).
$$

By similar arguments as in our current proof, we expect this assumption will be sufficient to bound $R_1$ for this modified algorithm. We expect that the bound on $R_2$ would be very similar as in the proof of Theorem 2, as this regret term corresponds to concentration of the cost. Similarly, the regret term $R_0$ can also be bounded the same as in the proof of Theorem 2, as this term corresponds to the warm-up period which still has length $\tilde{O}(T^{2/3})$. Therefore, we expect that the total regret of this modified algorithm can still be bounded by $\tilde{O}(T^{2/3})$.

### I.2.2 Theorem 1

We leave whether or not Theorem 1 generalizes to higher dimensions in all situations as an open question. However, we will briefly outline a setting in which we do expect the result to generalize. Suppose that $m = n$ and that $A^*$ and $B^*$ are invertible and diagonalizable.

Algorithm 3 for higher dimensions would require the same changes as in the previous sub-subsection, which means that $C_s^{\text{alg}} = C_{K_{\text{opt}}(\hat{\theta}_s, T_s, \Delta, d - \tilde{\Theta}_T(\epsilon_s))}^{\hat{\theta}_s}$. The main new result that would be necessary is an analogue of Lemma 21 for higher dimensions. Intuitively, the result of Lemma 21 holds because Algorithm 3 will a constant fraction of the time use the non-linear control $u_t^{\text{safeU}}$ which allows for faster learning. The analogue for higher dimensions is to show that the modified algorithm will a constant fraction of the time use a non-linear control. A difficulty in higher dimension is that it is not sufficient to just be non-linear along one dimension. Instead, there must be sufficient non-linearity in all $m$ dimensions. Therefore, the higher dimensional version of Assumption 9 requires that the noise distribution is sufficiently large relative to the constraints in all $m$ dimensions, which for example would be satisfied by the multivariate normal distribution with mean 0 and constant variance matrix. Under this assumption, the modified algorithm will a constant fraction of the time use controls $u_t$ that satisfy $\Delta_i(A^* x_t + B^* u_t) \geq d_i - O_T(\epsilon_s)$ for some $i \in [1 : k]$. Furthermore, if the noise is sufficiently large in all dimensions, then we expect that for every side of the boundary of the convex compact region (corresponding to $\Delta_i$ and $d_i$ for $i \in [1 : k]$), $x_t$ will at times be sufficiently far from that side and a point on that side will be the closest point to $x_t$. Because $A^*$ is invertible, the previous sentence will also hold for $A^* x_t$. Because $B^* u_t$ must bring the position back to within the safe region in expectation, for every side of the boundary we must have that $B^* u_t$ is large and perpendicular to that side. Because $B^*$ is invertible, this implies that the $u_t$ used to enforce safety will be sufficiently non-linear in all directions. We believe this would allow the algorithm to learn the matrix $B^*$ up to accuracy $O_T(1/\sqrt{t})$ at time $t$. Equipped with an analogue of Lemma 21, we expect that the rest of the proof will follow directly. If $m > n$ or $A^*$ and $B^*$ are not invertible, then showing that the non-linear controls $u_t$ are sufficient for learning every column of the matrix $B^*$ is more difficult. We leave the details of analyzing this case for future work.