

# Toward Understanding State Representation Learning in MuZero: A Case Study in Linear Quadratic Gaussian Control

**Yi Tian**

Massachusetts Institute of Technology

yitian@mit.edu

**Kaiqing Zhang**

University of Maryland, College Park

kaiqing@umd.edu

**Russ Tedrake**

Massachusetts Institute of Technology

russt@mit.edu

**Suvrit Sra**

Massachusetts Institute of Technology

suvrit@mit.edu

## Abstract

We study the problem of representation learning for control from partial and potentially high-dimensional observations. We approach this problem via *direct latent model learning*, where one directly learns a dynamical model in some latent state space by *predicting costs*. In particular, we establish *finite-sample guarantees* of finding a near-optimal representation function and a near-optimal controller using the directly learned latent model for infinite-horizon time-invariant Linear Quadratic Gaussian (LQG) control. A part of our approach to latent model learning closely resembles *MuZero*, a recent breakthrough in empirical reinforcement learning, in that it learns latent dynamics *implicitly* by predicting *cumulative costs*. A key technical contribution of this work is to prove persistency of excitation for a new stochastic process that arises from the analysis of quadratic regression in our approach.

## 1 Introduction

Control with a learned latent model is state-of-the-art in several reinforcement learning (RL) benchmarks, including board games, Atari games, and visuomotor control (Schrittwieser et al., 2020; Ye et al., 2021; Hafner et al., 2023). To better understand this modern machinery, we introduce it to a classical optimal control problem, namely Linear Quadratic Gaussian (LQG) control, and study its theoretical, finite-sample performance. Essential to this approach is the learning of two components: a *state representation* function that maps an observed history to some latent state, and a *latent model* that predicts the transition and cost in the latent state space. The latent model is usually a Markov decision process, using which we obtain a policy in the latent space or execute online planning.

What is the correct *objective* to optimize for learning a latent model? One popular choice is to learn a function that *reconstructs the observation* from the latent state (Hafner et al., 2019a,b,

2020, 2023). A latent model learned this way is *agnostic to control tasks* and retains all the information about the environment. This class of approaches can achieve satisfactory performance, but are prone to background distraction and control-irrelevant information (Fu et al., 2021). The second class of methods learn an *inverse model* that infers actions from latent states at different time steps (Pathak et al., 2017; Lamb et al., 2022). A latent model learned with this methodology is also task agnostic but extracts control-relevant information. In contrast, *task-relevant* representations can be learned by *predicting costs* in the control task (Oh et al., 2017; Zhang et al., 2020; Schrittwieser et al., 2020). The concept that a good latent state should be able to predict costs is intuitive, and the costs are directly relevant to optimal control. Hence, (Tian et al., 2022) refers to this class of methods as “direct latent model learning”, which is the focus of this work.

The direct latent model learning method of particular interest to us is that of MuZero (Schrittwieser et al., 2020). Announced by DeepMind in 2019, MuZero extends the line of works including AlphaGo (Silver et al., 2016), AlphaGo Zero (Silver et al., 2017), and AlphaZero (Silver et al., 2018) by not requiring knowledge of the game rules. MuZero matches the superhuman performance of AlphaZero in Go, shogi and chess, while outperforming model-free RL algorithms in Atari games. MuZero builds upon the powerful planning procedure of Monte Carlo Tree Search, with the major innovation being *learning a latent model*. The latent model replaces the rule-based simulator during planning, and avoids the burdensome planning in pixel space for Atari games.

MuZero is a milestone algorithm in representation learning for control. Intuitively, the algorithm design makes sense, but its complexity has so far inhibited a formal theoretical study. On the other hand, statistical learning theory for linear dynamical systems and control has evolved rapidly in recent years (Tsiamis et al., 2022); for partially observable linear dynamical systems, much of the work relies on learning *Markov parameters*, lacking a direct connection to the empirical methods used in practice for possibly nonlinear systems. In this work, we aim to bridge the two areas by studying provable MuZero-style latent model learning in LQG control.

The latent model learning of MuZero features three ingredients: 1) stacking frames, i.e., observations, as input to the representation function; 2) predicting costs, “optimal” values, and “optimal” actions from latent states; and 3) implicit learning of latent dynamics by predicting these quantities from latent states at future time steps. These are the defining characteristics of the MuZero-style algorithm that we shall consider. In MuZero, the “optimal” values and actions are found by the powerful online planning procedure. In this work, we simplify the setup by considering data collected using random actions, which are known to suffice for identifying a partially observable linear dynamical system (Oymak and Ozay, 2019). In this setup, the values become those associated with this trivial policy and we do not predict actions since they are random noises anyway. Note that although our study of the above ingredients is directly motivated by MuZero, previous empirical works have also explored them. For example, frame stacking has been a widely used technique to handle partial observability (Mnih et al., 2013, 2015); predicting values for learning a latent model has been studied in (Oh et al., 2017), which also learns the latent state transition implicitly.

Closely related to our work, (Tian et al., 2022) also considers provable direct latent model learning in LQG, but for the finite-horizon time-varying setting. Our work builds upon it

and complements it in two ways: 1) we extend their algorithm to the time-invariant setting with a *stationary* representation function and latent model, which is closer to what has been deployed in practice; 2) we present and analyze a new, MuZero-style latent model learning algorithm. Both 1) and 2) introduce new technical challenges to be addressed. We summarize our contributions as follows.

- We show that two direct latent model learning methods provably solve infinite-horizon time-invariant LQG control by establishing finite-sample guarantees. Both methods only need a single trajectory; one resembles the method in (Tian et al., 2022), and the other resembles MuZero.
- By analyzing the MuZero-style algorithm, we notice the potential issue of coordinate misalignment; that is, costs can be invariant to certain transformations of the latent states, and implicit dynamics learning by predicting *one-step* transition may not recover the latent state coordinates consistently. This insight suggests the need of predicting *multi-step* transition or other coordinate alignment procedures in implicit dynamics learning.
- Technically, we overcome the difficulty of having *dependent* data samples in a single trajectory for latent model learning, by proving a new result about the persistency of excitation for a stochastic process that arises from the analysis of the quadratic regression subroutine in both of our methods.

**Notation.** Random vectors are denoted by lowercase letters; sometimes they also denote their realized values. Uppercase letters denote matrices, some of which can be random. Let  $a \wedge b$  denote the minimum between scalars  $a$  and  $b$ .  $\mathbf{1}$  denotes either the scalar one or a vector consisting of all ones;  $I$  denotes an identity matrix. The dimension, when emphasized, is specified in subscripts, e.g.,  $\mathbf{1}_d, I_d$ . Given vector  $v \in \mathbb{R}^d$ , let  $\|v\|$  denote its  $\ell_2$  norm and  $\|v\|_P := (v^\top P v)^{1/2}$  for positive semidefinite  $P \in \mathbb{R}^{d \times d}$ . Given symmetric matrices  $P$  and  $Q$ ,  $P \succ Q$  or  $Q \prec P$  means  $P - Q$  is positive definite, and  $P \succcurlyeq Q$  or  $Q \preccurlyeq P$  means  $P - Q$  is positive semidefinite. Semicolon “;” denotes stacking vectors or matrices vertically. For a collection of  $d$ -dimensional vectors  $(v_t)_{t=i}^j$ , let  $v_{i:j} := [v_i; v_{i+1}; \dots; v_j] \in \mathbb{R}^{d(j-i+1)}$  denote the concatenation along the column. For random variable  $x$ , let  $\|x\|_{\psi_\beta}$  denote its  $\beta$ -sub-Weibull norm, a special case of Orlicz norms (Zhang and Wei, 2022), with  $\beta = 1, 2$  corresponding to subexponential and sub-Gaussian norms. For matrix  $A$ , let  $\sigma_{\min}(A)$ ,  $\|A\|_2$ ,  $\|A\|_F$ , and  $\|A\|_*$  denote its minimum eigenvalue, minimum singular value, operator norm (induced by vector  $\ell_2$  norms), Frobenius norm, and nuclear norm, respectively.  $\langle \cdot, \cdot \rangle_F$  denotes the Frobenius inner product between matrices. For square matrix  $A$ , let  $\lambda_{\min}(A)$  be its minimum eigenvalue and  $\rho(A)$  be its spectral radius. Define  $\alpha(A) := \sup_{k \geq 0} \|A^k\|_2 \rho(A)^{-k}$ . Let  $\text{svec}(\cdot)$  denote the operator of flattening a symmetric matrix by stacking its columns; it does not repeat the off-diagonal elements, but scales them by  $\sqrt{2}$  (Schacke, 2004). We adopt the standard use of  $\mathcal{O}(\cdot), \Omega(\cdot), \Theta(\cdot)$ , where the hidden constants are dimension-free but may depend on system parameters.

## 2 Problem setup

A partially observable linear time-invariant (LTI) dynamical system is described by

$$x_{t+1} = A^* x_t + B^* u_t + w_t, \quad y_t = C^* x_t + v_t, \quad (2.1)$$

with state  $x_t \in \mathbb{R}^{d_x}$ , observation  $y_t \in \mathbb{R}^{d_y}$ , and control  $u_t \in \mathbb{R}^{d_u}$  for all  $t \geq 0$ . Process noises  $(w_t)_{t \geq 0}$  and observation noises  $(v_t)_{t \geq 0}$  are i.i.d. zero-mean Gaussian random vectors with covariance matrices  $\Sigma_w$  and  $\Sigma_v$ , respectively, and the two sequences are mutually independent. Let initial state  $x_0$  be sampled from  $\mathcal{N}(0, \Sigma_0)$ . The quadratic cost function is given by

$$c(x, u) = \|x\|_{Q^*}^2 + \|u\|_{R^*}^2, \quad (2.2)$$

where  $Q^* \succcurlyeq 0$  and  $R^* \succ 0$ .

A policy/controller  $\pi$  determines an action/control input  $u_t$  at time step  $t$  based on the history  $[y_{0:t}; u_{0:(t-1)}]$  up to this time step. For  $t \geq 0$ ,  $c_t := c(x_t, u_t)$  denotes the cost at time step  $t$ . Given a policy  $\pi$ , let

$$J^\pi := \limsup_{T \rightarrow \infty} \mathbb{E} \left[ \frac{1}{T} \sum_{t=0}^{T-1} c_t \right] \quad (2.3)$$

denote the time-average expected cost. The objective of LQG control is to find a policy  $\pi$  such that  $J^\pi$  is minimized.

In the fully observable setting, known as the linear quadratic regulator (LQR),  $y_t = x_t$ . A linear controller with feedback gain  $K \in \mathbb{R}^{d_u \times d_x}$  determines action  $u_t = Kx_t$  at time step  $t$ . Let  $J^K(A^*, B^*, Q^*, R^*)$  denote the time-average expected cost (2.3) in the LQR problem  $(A^*, B^*, Q^*, R^*)$  under feedback gain  $K$  and define  $J^*(A^*, B^*, Q^*, R^*) := \min_K J^K(A^*, B^*, Q^*, R^*)$ .

We make the following standard assumptions.

**Assumption 1.** System dynamics (2.1) and cost (2.2) satisfy:

1. The system is stable, that is,  $\rho(A^*) < 1$ .
2.  $(A^*, B^*)$  is  $\nu$ -controllable for some  $\nu > 0$ , that is, the controllability matrix

$$\Phi_c(A^*, B^*) := [B^*, A^* B^*, \dots, (A^*)^{d_x-1} B^*]$$

has rank  $d_x$  and  $\sigma_{\min}(\Phi_c(A^*, B^*)) \geq \nu$ .

3.  $(A^*, C^*)$  is  $\omega$ -observable for some  $\omega > 0$ , that is, the observability matrix

$$\Phi_o(A^*, C^*) := [C^*; C^* A^*; \dots; C^* (A^*)^{d_x-1}]$$

has rank  $d_x$  and  $\sigma_{\min}(\Phi_o(A^*, C^*)) \geq \omega$ .

4.  $(A^*, \Sigma_w^{1/2})$  is  $\eta$ -controllable for some  $\eta > 0$ .
5.  $(A^*, (Q^*)^{1/2})$  is  $\mu$ -observable for some  $\mu > 0$ .

6.  $\Sigma_v \succcurlyeq \sigma_v^2 I$  for some  $\sigma_v > 0$ ; this can always be achieved by inserting Gaussian noises with full-rank covariance matrices to the observations.
7.  $R^* \succcurlyeq r^2 I$  for some  $r > 0$ .
8. The operator norms of  $A^*, B^*, C^*, Q^*, R^*, \Sigma_w, \Sigma_v, \Sigma_0$  are  $\mathcal{O}(1)$  and the singular value lower bounds  $\nu, \omega, \nu, \eta, \sigma_v, r$  are  $\Omega(1)$ .

If the system parameters  $(A^*, B^*, C^*, Q^*, R^*, \Sigma_w, \Sigma_v)$  are known, the optimal policy is obtained by combining the Kalman filter

$$z_{t+1}^* = A^* z_t^* + B^* u_t + L^* (y_{t+1} - C^* (A^* z_t^* + B^* u_t)) \quad (2.4)$$

with the optimal feedback gain  $K^*$  of the linear quadratic regulator (LQR) such that  $u_t = K^* z_t^*$ , where  $L^*$  is the Kalman gain, and at the initial time step, we can set, e.g.,  $z_0^* = L^* y_0$ . This fact is known as the *separation principle*, and the Kalman gain and optimal feedback gain are given by

$$L^* = S^* (C^*)^\top (C^* S^* (C^*)^\top + \Sigma_v)^{-1}, \quad (2.5)$$

$$K^* = -((B^*)^\top P^* B^* + R)^{-1} (B^*)^\top P^* A^*, \quad (2.6)$$

where  $S^*$  and  $P^*$  are determined by their respective discrete-time algebraic Riccati equations (DAREs):

$$S^* = A^* (S^* - S^* (C^*)^\top (C^* S^* (C^*)^\top + \Sigma_v)^{-1} C^* S^*) (A^*)^\top + \Sigma_w, \quad (2.7)$$

$$P^* = (A^*)^\top (P^* - P^* B^* ((B^*)^\top P^* B^* + R)^{-1} (B^*)^\top P^*) A^* + Q^*. \quad (2.8)$$

Assumptions 1.2 to 1.7 guarantee the existence and uniqueness of positive definite solutions  $S^*$  and  $P^*$ ; Assumption 1.8 further guarantees that their operator norms are  $\mathcal{O}(1)$  and minimum singular values are  $\Omega(1)$ .

We consider the data-driven control setting, where the LQG model  $(A^*, B^*, C^*, Q^*, \Sigma_w, \Sigma_v)$  is unknown. For simplicity, we assume  $R^*$  is known, though our approaches can be readily extended to the case where it is unknown by learning it from predicting costs.

## 2.1 Latent model of LQG

The stationary Kalman filter (2.4) asymptotically produces the optimal *state estimation* in the sense of minimum mean squared errors. With a finite horizon, however, the optimal state estimator is time-varying, given by

$$z_{t+1}^* = A^* z_t^* + B^* u_t + L_{t+1}^* (y_{t+1} - C^* (A^* z_t^* + B^* u_t)), \quad (2.9)$$

where  $L_t^*$  is the time-varying Kalman gain, converging to  $L^*$  as  $t \rightarrow \infty$ . This convergence is equivalent to that of error covariance matrix  $\mathbb{E}[(x_t - z_t^*)(x_t - z_t^*)^\top]$ , which happens exponentially fast (Komaroff, 1994). Hence, for simplicity, we assume this error covariance matrix is stationary at the initial time step by the choice of  $z_0^*$  so that  $L_t^* = L^*$  for  $t \geq 1$ ; this assumption is common in the literature (Lale et al., 2020, 2021; Jadbabaie et al., 2021). The *innovation* term  $i_{t+1} := y_{t+1} - C^* (A^* z_t^* + B^* u_t)$  is independent of the history  $(y_0, u_0, y_1, u_1, \dots, y_{t+1})$  and  $(i_t)_{t \geq 1}$

are mutually independent. The following proposition taken from (Tian et al., 2022, Proposition 1) represents the system in terms of the state estimates obtained by the Kalman filter, which we refer to as the *latent model*.

**Proposition 1.** *Let  $(z_t^*)_{t \geq 1}$  be state estimates given by the time-varying Kalman filter. Then, for  $t \geq 0$ ,*

$$z_{t+1}^* = A^* z_t^* + B^* u_t + L^* i_{t+1},$$

where  $L^* i_{t+1}$  is independent of  $z_t^*$  and  $u_t$ , i.e., the state estimates follow the same linear dynamics with noises  $L^* i_{t+1}$ . The cost at step  $t$  can be reformulated as functions of the state estimates by

$$c_t = \|z_t^*\|_{Q^*}^2 + \|u_t\|_{R^*}^2 + b^* + \gamma_t + \eta_t,$$

where  $b^* > 0$ , and  $\gamma_t = \|z_t^* - x_t\|_{Q^*}^2 - b^*$ ,  $\eta_t = \langle z_t^*, x_t - z_t^* \rangle_{Q^*}$  are both zero-mean subexponential random variables. Moreover,  $b^* = \mathcal{O}(1)$  and  $\|\gamma_t\|_{\psi_1} = \mathcal{O}(d_x^{1/2})$ ; if control  $u_t \sim \mathcal{N}(0, \sigma_u^2 I)$  for  $t \geq 0$ , then we have  $\|\eta_t\|_{\psi_1} = \mathcal{O}(d_x^{1/2})$ .

Proposition 1 shows that the dynamics of the state estimates computed by the time-varying Kalman filter is the same as the original system up to noises; the costs are also the same, up to constants and noises. Hence, a latent model can be parameterized by  $(A, B, Q, R^*)$ , with the constant  $b^*$  and noises neglected due to their irrelevance to planning. A stationary latent policy is a linear controller  $u_t = Kz_t$  on latent state  $z_t$ , parameterized by feedback gain  $K \in \mathbb{R}^{d_u \times d_x}$ .

The latent model enables us to find a good latent policy. To learn such a latent model and to deploy a latent policy in the original partially observable system, we need a representation function. Let  $\bar{A}^* := (I - L^* C^*) A^*$  and  $\bar{B}^* := (I - L^* C^*) B^*$ . Then, the Kalman filter can be written as  $z_{t+1}^* = \bar{A}^* z_t^* + \bar{B}^* u_t + L^* y_{t+1}$ . For  $t \geq 0$ , unrolling the recursion gives

$$\begin{aligned} z_t^* &= \bar{A}^* (\bar{A}^* z_{t-2}^* + \bar{B}^* u_{t-2} + L^* y_{t-1}) + \bar{B}^* u_{t-1} + L^* y_t \\ &= [(\bar{A}^*)^{t-1} L^*, \dots, L^*] y_{1:t} + [(\bar{A}^*)^{t-1} \bar{B}^*, \dots, \bar{B}^*] u_{0:(t-1)} + (\bar{A}^*)^t z_0^* \\ &=: M_t^* [y_{1:t}; u_{0:(t-1)}; z_0^*], \end{aligned}$$

where  $M_t^* \in \mathbb{R}^{d_x \times (td_y + td_u + d_x)}$ . This means the representation function can be parameterized as linear mappings for full histories (with  $y_0$  replaced by  $z_0^*$ ). Despite the simplicity, the input dimension of the function grows linearly in time, making it intractable to estimate the state using the full history for large  $t$ ; nor is it necessary, since the impact of old data decreases exponentially. Under Assumption 1,  $\rho(\bar{A}^*) < 1$  (Bertsekas, 2012, Appendix E.4). With an  $H$ -step truncated history, the state estimate can be written as

$$\begin{aligned} z_t^* &= [(\bar{A}^*)^{H-1} L^*, \dots, L^*] y_{(t-H+1):t} + [(\bar{A}^*)^{H-1} \bar{B}^*, \dots, \bar{B}^*] u_{(t-H):(t-1)} + \delta_t \\ &=: M^* [y_{(t-H+1):t}; u_{(t-H):(t-1)}] + \delta_t, \end{aligned}$$

where  $\delta_t = (\bar{A}^*)^H z_{t-H}^*$ , whose impact decays exponentially in  $H$  and can be neglected for sufficiently large  $H$ , since  $z_{t-H}^*$  converges to a stationary distribution and its norm is bounded with high probability. Hence, the representation function that we aim to recover is  $M^* \in$



$\mathbb{R}^{d_x \times H(d_y + d_u)}$ , which takes as input the  $H$ -step history  $h_t = [y_{(t-H+1):t}; u_{(t-H):(t-1)}]$ . Henceforth, we let  $d_h := H(d_y + d_u)$ . Then, a representation function is parameterized by matrix  $M \in \mathbb{R}^{d_x \times d_h}$ .

Overall, a policy is a combination of a representation function  $M$  and a feedback gain  $K$  in the latent model, denoted by  $\pi = (M, K)$ . Learning to solve LQG control in this framework can thus be achieved by: 1) learning representation function  $M$ ; 2) extracting latent model  $(A, B, Q, R^*)$ ; and 3) finding the optimal  $K$  by planning in the latent model. Next, we introduce our approach following this pipeline.

### 3 Method

In practice, latent model learning methods collect trajectories by interacting with the system online using some policy; the trajectories are used to improve the learned latent model, which in turn improves the policy. In LQG control, it is known that the simple setup allows us to learn a good latent model from a single trajectory, collected using zero-mean Gaussian inputs; see e.g., (Oymak and Ozay, 2019). This is also how we assume the data are collected. We note that our results also apply to data from multiple independent trajectories using the same zero-mean Gaussian inputs.

In direct latent model learning, state representations are learned by predicting costs. To learn the transition function in the latent model, two approaches are explored in the literature. The first approach explicitly minimizes transition prediction errors (Subramanian et al., 2020; Hafner et al., 2019a). Algorithmically, the overall loss is a combination of cost and transition prediction errors. The second approach, which MuZero takes, learns transition *implicitly*, by minimizing *cost prediction errors at future states* generated from the transition function (Oh et al., 2017; Schrittwieser et al., 2020). Algorithmically, the overall loss aggregates the cost prediction errors across multiple time steps. In both approaches, the coupling of different terms in the loss makes finite-sample analysis difficult. As observed in (Tian et al., 2022), the structure of LQG allows us to learn the representation function independently of learning the transition function. This allows us to formulate both approaches under the same direct latent model learning framework (Algorithm 1).

Algorithm 1 consists of three main steps. Lines 3 to 5 correspond to cost-driven representation learning. Lines 6 to 8 correspond to latent model learning, where the system dynamics can be identified either explicitly, by ordinary least squares (SysId), or implicitly, by future cost prediction (CoSysId, Algorithm 2). Line 8 corresponds to latent policy optimization; in LQG this amounts to solving a DARE. Below we elaborate on cost-driven representation learning, SysId, and CoSysId in order.

#### 3.1 Cost-driven representation learning

The procedure of cost-driven representation learning is almost identical to that in (Tian et al., 2022). The main idea is to perform quadratic regression (3.2) to the  $d_x$ -step cumulative costs; these correspond to the value prediction in MuZero. By the  $\mu$ -observability of  $(A^*, (Q^*)^{1/2})$

---

**Algorithm 1** Direct latent model learning for LQG control

---

- 1: **Input:** length  $T$ , history length  $H$ , noise magnitude  $\sigma_u$
- 2: Collect a trajectories of length  $T + H$  using  $u_t \sim \mathcal{N}(0, \sigma_u^2 I)$ , for  $t \geq 0$ , to obtain  $\mathcal{D}_{\text{raw}}$ :

$$(y_0, u_0, c_0, y_1, u_1, c_1, \dots, y_{T+H}) \quad (3.1)$$

- 3: Estimate the state representation function and cost constants by solving

$$\hat{N}, \hat{b}_0 \in \underset{N=N^\top, b_0}{\operatorname{argmin}} \sum_{t=H}^{T+H-1} (\|h_t\|_N^2 + b_0 - \bar{c}_t)^2, \quad (3.2)$$

where  $\bar{c}_t := \sum_{\tau=t}^{t+d_x-1} (c_\tau - \|u_\tau\|_{R^*}^2)$

- 4: Find  $\hat{M} \in \underset{M \in \mathbb{R}^{d_x \times H(d_y+d_u)}}{\operatorname{argmin}} \|M^\top M - \hat{N}\|_F$
- 5: Compute  $\hat{z}_t = \hat{M}[y_{(t-H+1):t}; u_{(t-H):(t-1)}]$  for all  $t \geq H$ , so that the data are converted to  $\mathcal{D}_{\text{state}}$ :

$$(\hat{z}_H, u_H, c_H, \dots, \hat{z}_{T+H-1}, u_{T+H-1}, c_{T+H-1}, \hat{z}_{T+H})$$

- 6: Run SysID (3.3) or CoSysID (Algorithm 2) to obtain dynamics matrices  $(\hat{A}, \hat{B})$
- 7: Estimate the cost function by solving

$$\tilde{Q}, \hat{b} \in \underset{Q=Q^\top, b}{\operatorname{argmin}} \sum_{t=H}^{T+H-1} (\|\hat{z}_t\|_Q^2 + b - c_t)^2,$$

- 8: Truncate negative eigenvalues of  $\tilde{Q}$  to zero to obtain  $\hat{Q} \succcurlyeq 0$
  - 9: Find feedback gain  $\hat{K}$  from  $(\hat{A}, \hat{B}, \hat{Q}, R^*)$  by DARE (2.8) and (2.6)
  - 10: **Return:** policy  $\hat{\pi} = (\hat{M}, \hat{K})$
- 

(Assumption 1.5), the cost observability Gramian

$$\bar{Q}^* := \sum_{t=0}^{d_x-1} ((A^*)^t)^\top Q^* (A^*)^t \succcurlyeq \mu^2 I.$$

Under zero control and zero noise, starting from  $x$ , the  $d_x$ -step cumulative cost is precisely  $\|x\|_{\bar{Q}^*}^2$ . Hence,  $\hat{N}$  estimates  $N^* = (M^*)^\top \bar{Q}^* M^*$ ; up to an orthonormal transformation,  $\hat{M}$  recovers  $M^{*'} := (\bar{Q}^*)^{1/2} M^*$ , the representation function under an equivalent parameterization, termed as the *normalized parameterization* in (Tian et al., 2022), where

$$A^{*'} = (\bar{Q}^*)^{1/2} A^* (\bar{Q}^*)^{-1/2}, \quad B^{*'} = (\bar{Q}^*)^{1/2} B, \quad C^{*'} = C^* (\bar{Q}^*)^{-1/2}, \\ w_t' = (\bar{Q}^*)^{1/2} w_t, \quad Q^{*'} = (\bar{Q}^*)^{-1/2} Q^* (\bar{Q}^*)^{-1/2}.$$

Due to the following proposition, the algorithm does not need to know the dimension  $d_x$  of the latent model; it can discover  $d_x$  from the eigenvalues of  $\hat{N}$ .

**Proposition 2.** Under i.i.d. control inputs  $u_t \sim \mathcal{N}(0, \sigma_u^2 I)$  for  $t \geq 0$ ,  $\lambda_{\min}(\operatorname{Cov}(z_t^*)) = \Omega(v^2)$  for  $t \geq d_x$ , where  $v$  is defined in Assumption 1.3. Recall that for a square matrix  $A$ , we define  $\alpha(A) := \sup_{k \geq 0} \|A^k\|_2 \rho(A)^{-k}$ . As long as  $H \geq \frac{\log(a\alpha(\bar{A}^*))}{\log(\rho(\bar{A}^*)^{-1})}$  for some dimension-free constant  $a > 0$ ,  $M^*$  has rank  $d_x$  and  $\sigma_{\min}(M^*) \geq \Omega(vH^{-1/2})$ .



*Proof.* For  $t \geq d_x$ , unrolling the Kalman filter gives

$$\begin{aligned} z_t^* &= A^* z_{t-1}^* + B^* u_{t-1} + L^* i_t \\ &= A^* (A^* z_{t-2}^* + B^* u_{t-2} + L^* i_{t-1}) + L^* i_t \\ &= [B^*, \dots, (A^*)^{d_x-1} B^*] [u_{t-1}; \dots; u_{t-d_x}] + (A^*)^{d_x} z_{t-d_x}^* + [L^*, \dots, (A^*)^{d_x-1} L^*] [i_t; \dots; i_{t-d_x+1}], \end{aligned}$$

where  $(u_\tau)_{\tau=t-d_x}^{t-1}$ ,  $z_{t-d_x}^*$  and  $(i_\tau)_{\tau=t-d_x+1}^t$  are independent. For  $H \geq d_x$ , the matrix multiplied by  $[u_{t-1}; \dots; u_{t-d_x}]$  is precisely the controllability matrix  $\Phi_c(A^*, B^*)$ . Then,

$$\begin{aligned} \text{Cov}(z_t^*) &= \mathbb{E}[z_t^* (z_t^*)^\top] \succcurlyeq \Phi_c(A^*, B^*) \mathbb{E}[[u_{t-1}; \dots; u_{t-d_x}][u_{t-1}; \dots; u_{t-d_x}]^\top] \Phi_c^\top(A^*, B^*) \\ &= \sigma_u^2 \Phi_c(A^*, B^*) \Phi_c^\top(A^*, B^*). \end{aligned}$$

By the  $\nu$ -controllability of  $(A^*, B^*)$ ,  $\text{Cov}(z_t^*)$  is full-rank and  $\lambda_{\min}(\text{Cov}(z_t^*)) \geq \sigma_u^2 \nu^2$ . Since  $z_t^* = M^* h_t + \delta_t$ , we have

$$\text{Cov}(M^* h_t) = \text{Cov}(z_t^* - \delta_t) = \text{Cov}(z_t^*) + \text{Cov}(\delta_t) - \text{Cov}(z_t^*, \delta_t) - \text{Cov}(\delta_t, z_t^*).$$

By Lemma 3,

$$\|\text{Cov}(z_t^*, \delta_t)\|_2 = \|\text{Cov}(\delta_t, z_t^*)\|_2 \leq \|\text{Cov}(z_t^*)^{1/2}\|_2 \|\text{Cov}(\delta_t)^{1/2}\|_2.$$

Hence, by Weyl's inequality,

$$\lambda_{\min}(\text{Cov}(M^* h_t)) \geq \lambda_{\min}(\text{Cov}(z_t^*)) - 2\|\text{Cov}(z_t^*)^{1/2}\|_2 \|\text{Cov}(\delta_t)^{1/2}\|_2.$$

Since  $\|\text{Cov}(z_t^*)\|_2 = \mathcal{O}(1)$  due to the stability of  $A^*$  and  $\delta_t = (\bar{A}^*)^H z_{t-H}^*$ , there exists some dimension-free constant  $a > 0$  such that as long as  $H \geq \frac{\log(aa(\bar{A}^*))}{\log(\rho(\bar{A}^*)^{-1})}$ ,

$$\lambda_{\min}(\text{Cov}(M^* h_t)) \geq \sigma_u^2 \nu^2 / 2.$$

On the other hand,

$$\mathbb{E}[M^* h_t h_t^\top (M^*)^\top] \preccurlyeq \|\mathbb{E}[h_t h_t^\top]\|_2 M^* (M^*)^\top.$$

Since  $h_t = [y_{(t-H+1):t}; u_{(t-H):(t-1)}]$  and  $(\text{Cov}(y_t))_{t \geq 0}$ ,  $(\text{Cov}(u_t))_{t \geq 0}$  have  $\mathcal{O}(1)$  operator norms, by Lemma 4,  $\|\text{Cov}(h_t)\|_2 = \|\mathbb{E}[h_t h_t^\top]\|_2 = \mathcal{O}(H)$ . Hence,

$$0 < \sigma_u^2 \nu^2 / 2 \leq \lambda_{\min}(\text{Cov}(M^* h_t)) = \mathcal{O}(H) \sigma_{d_x}^2(M^*).$$

Since  $M^* \in \mathbb{R}^{d_x \times d_h}$ , this implies that  $\text{rank}(M^*) = d_x$  and  $\sigma_{\min}(M^*) = \Omega(\nu H^{-1/2})$ .  $\square$

Proposition 2 is an adaption of (Tian et al., 2022, Proposition 2) to the infinite-horizon LTI setting. Necessarily, this implies that by our choice of  $H$ ,  $d_h = H(d_y + d_u) \geq d_x$ . Moreover, since  $\bar{Q}^* \succcurlyeq \mu^2 I$ ,  $N^* = (M^*)^\top \bar{Q}^* M^*$  is a  $d_h \times d_h$  matrix with rank  $d_x$ , and  $\lambda_{\min}^+(N^*) \geq \lambda_{\min}(\bar{Q}^*) \lambda_{\min}^2(M^*) = \Omega(\mu^2 \nu^2 H^{-1})$ . Hence, if  $\hat{N}$  is sufficiently close to  $N^*$ , by setting an appropriate threshold on the eigenvalues of  $\hat{N}$ , the dimension of the latent model equals the number of eigenvalues above it.

To find an approximate factorization of  $\hat{N}$ , let  $\hat{N} = U\Lambda U^\top$  be its eigenvalue decomposition, where the diagonal elements of  $\Lambda$  are listed in a descending order, and  $U$  is an orthonormal matrix. Let  $\Lambda_{d_x}$  be the left-top block of  $\Lambda$  and  $U_{d_x}$  be the left  $d_x$  columns of  $U$ . By the Eckart-Young-Mirsky theorem,  $\hat{M} = \max(\Lambda_{d_x}, 0)^{1/2} U_{d_x}^\top$ , where “max” applies elementwise, is the solution to Line 4 of Algorithm 1, that is, the best approximate factorization of  $\hat{N}$  among  $d_x \times d_h$  matrices in terms of the Frobenius norm approximation error.

In the next two subsections, we move on to discuss learning latent dynamics, including the explicit approach SysID and the implicit approach CoSysID.

### 3.2 Explicit learning of system dynamics

Explicit learning of the system dynamics simply minimizes the transition prediction error in the latent space (Subramanian et al., 2020), or more generally, the statistical distances between the predicted and estimated distributions of the next latent state, like the KL divergence (Hafner et al., 2019a). In linear systems, it suffices to use the ordinary least squares as the SysID procedure, that is, to solve

$$(\hat{A}, \hat{B}) \in \underset{A, B}{\operatorname{argmin}} \sum_{t=H}^{T+H-1} \|A\hat{z}_t + Bu_t - \hat{z}_{t+1}\|^2. \quad (3.3)$$

In this linear regression, if  $(\hat{z}_t)_{t \geq H}$  are the optimal state estimates  $(z_t^*)_{t \geq H}$  (2.9), then (Simchowitz et al., 2018) has shown finite-sample guarantees for  $(\hat{A}, \hat{B})$ . Here,  $\hat{z}_t$  contains errors resulting from the representation function  $\hat{M}$  and the residual error  $\delta_t$ , but as long as  $T$  and  $H$  are large enough, SysID still has a finite-sample guarantee, as will be shown in Lemma 8. We refer to the algorithm that instantiates Algorithm 1 with SysID as CoREL (Cost-driven state Representation Learning). As the time-varying counterpart in (Tian et al., 2022), it provably solves learning for LQG control, as will be shown in Theorem 2.

### 3.3 Implicit learning of system dynamics (MuZero-style)

An important ingredient of latent model learning in MuZero (Schrittwieser et al., 2020) is to *implicitly* learn the transition function by minimizing the cost prediction error at *future latent states* generated from the transition function. Let  $z_t = Mh_t$  denote the latent state given by representation function  $M$  at step  $t$ . Let  $z_{t,0} = z_t$  and  $z_{t,i} = Az_{t,i-1} + Bu_{t+i-1}$  for  $i \geq 1$  be the future latent state predicted by dynamics  $(A, B)$  from  $z_t$  after  $i$  steps of transition. For a trajectory of length  $T + H$  like (3.1), the loss that considers  $\ell$  steps into the future is given by

$$\sum_{t=H}^{T+H-K-1} \sum_{i=0}^{\ell} (\|z_{t,i}\|_Q^2 + \|u_t\|_{R^*}^2 + b - c_t)^2.$$

This loss involves powers of  $A$  up to  $A^\ell$ ; with the squared norm, the powers double, making the minimization over  $A$  hard to solve and analyze for  $\ell \geq 2$ . In LQG control, our finding is that it suffices to take  $\ell = 1$ . As mentioned in §1, MuZero also predicts optimal values and optimal actions; in LQG, to handle  $Q^* \neq 0$ , like cost-driven representation learning (see §3.1),

---

**Algorithm 2** CoSysID: Cost-driven system identification
 

---

- 1: **Input:** data  $\mathcal{D}_{\text{raw}}$ , representation function  $\hat{M}$
- 2: Estimate the system dynamics by

$$\hat{N}_1, \hat{b}_1 \in \underset{N_1=N_1^\top, b_1}{\operatorname{argmin}} \sum_{t=H}^{T+H-1} (\| [h_t; u_t] \|_{N_1}^2 + b_1 - \bar{c}_{t+1})^2 \quad (3.6)$$

- 3: Find  $\hat{M}_1 \in \underset{M_1 \in \mathbb{R}^{d_x \times (Hd_y + (H+1)d_u)}}{\operatorname{argmin}} \| M_1^\top M_1 - \hat{N}_1 \|_F$
- 4: Split  $\hat{M}_1$  to  $[\tilde{M}, \tilde{B}]$  after column  $H(d_y + d_u)$  and set  $\tilde{A} = \tilde{M}\hat{M}^\dagger$ .
- 5: Find alignment matrix  $\hat{S}_0$  by

$$\hat{S}_0 \in \underset{S_0 \in \mathbb{R}^{d_x \times d_x}}{\operatorname{argmin}} \sum_{t=H}^{T+H-1} \| S_0 \hat{M}_1 [h_t; u_t] - \hat{z}_{t+1} \|^2$$

- 6: **Return:** system dynamics estimate  $(\hat{A}, \hat{B}) = (\hat{S}_0 \tilde{A}, \hat{S}_0 \tilde{B})$
- 

we adopt the *cumulative costs* and use the normalized parameterization. Thus, the optimization problem we aim to solve is given by

$$\min_{M, A, B, b} \sum_{t=H}^{T+H-1} ((\| M h_t \|^2 + b - \bar{c}_t)^2 + (\| A M h_t + B u_t \|^2 + b - \bar{c}_{t+1})^2). \quad (3.4)$$

To convexify the optimization problem (3.4), we define  $N := M^\top M$  and  $N' := [A M, B]^\top [A M, B]$ . Then, (3.4) becomes

$$\min_{N, N_1, b} \sum_{t=H}^{T+H-1} ((\| h_t \|_N^2 + b - \bar{c}_t)^2 + (\| [h_t; u_t] \|_{N_1}^2 + b - \bar{c}_{t+1})^2). \quad (3.5)$$

This minimization problem is convex in  $N$ ,  $N_1$  and  $b$ , and has a closed-form solution; essentially, it consists of two linear regression problems coupled by  $b$ . Since constant  $b$  is merely a term accounting for the estimation error and not part of the representation function, we can decouple the two regression problems by allowing  $b$  to take different values in them. This further simplifies the algorithm: the first regression problem is exactly cost-driven representation learning (§3.1), and the second is cost-driven system identification (CoSysID, Algorithm 2). The algorithm that instantiates Algorithm 1 with CoSysID is called CoREDyL (Cost-driven state Representation and Dynamic Learning). Like CoREL, this MuZero-style latent model learning method provably solves LQG control, as we will show in Theorem 2.

CoSysID has similar steps to cost-driven representation learning (§3.1), except that in Line 5, it requires fitting a matrix  $\hat{S}_0$ . This is because the approximate factorization steps recover  $M^*$  and  $M_1^*$  up to orthonormal transformations, but there is no guarantee for the two orthonormal matrices to be the same; we need to fit  $\hat{S}_0$  to align their coordinates. We note that although CoSysID needs the output  $\hat{M}$  from cost-driven representation learning, the two quadratic regressions (3.2) and (3.6) are not coupled and can be solved in parallel.

## 4 Theoretical guarantees

The following Theorem 1 shows that both CoReL and CoReDyL are guaranteed to solve unknown LQG control with a finite number of samples.

**Theorem 1.** *Given an unknown LQG problem satisfying Assumption 1, let  $M^{*f}$  and  $(A^{*f}, B^{*f}, Q^{*f}, R^*)$  be the optimal state representation function and the true system parameters under the normalized parameterization. For a given  $p \in (0, 1)$ , if we run CoReL (Algorithm 1 with (3.3)) or CoReDyL (Algorithm 1 with Algorithm 2) for  $T \geq \text{poly}(d_x, d_y, d_u, \log(T/p))$ ,  $H = \Omega(\log(H^2(d_y + d_u)T \log(T/p)))$ , and  $\sigma_u = \Theta(1)$ , then there exists an orthonormal matrix  $S \in \mathbb{R}^{d_x \times d_x}$ , such that with probability at least  $1 - p$ , the representation function  $\hat{M}$  satisfies*

$$\|\hat{M} - SM^{*f}\|_2 = \mathcal{O}(\text{poly}(H, d_x, d_u, d_y, \log(T/p))T^{-1/2}),$$

and the feedback gain  $\hat{K}$  satisfies

$$\begin{aligned} & J^{\hat{K}}(SA^{*f}S^\top, SB^{*f}, SQ^{*f}S^\top, R^*) - J^*(SA^{*f}S^\top, SB^{*f}, SQ^{*f}S^\top, R^*) \\ &= \mathcal{O}(\text{poly}(H, d_x, d_u, d_y, \log(T/p))T^{-1}). \end{aligned}$$

We defer the proof of Theorem 1 to §F. Compared with common system identification methods based on learning Markov parameters (Oymak and Ozay, 2019; Simchowitz et al., 2019), the error bounds of the system parameters produced by CoReDyL (or CoReL) have the same dependence on  $T$ , but worse dependence on system dimensions. Moreover, to establish persistency of excitation, CoReDyL (or CoReL) requires a larger burn-in period. These relative sample inefficiencies are the price we pay for direct latent model learning, which is only supervised by *scalar-valued* costs that are quadratic in the history, instead of *vector-valued* observations that are linear in the history. Hence, we have to address the more challenging problem of *quadratic regression*, which lifts the dimension of the optimization problem. On the other hand, direct latent model learning avoids learning the reconstruction function  $C^*$  and can learn task-relevant representations in more complex settings, as demonstrated by empirical studies.

### 4.1 Persistency of excitation

Central to the analysis of CoReL and CoReDyL is the finite-sample characterization of the *quadratic regression* problem. To solve (3.2), notice that

$$\|h_t\|_N^2 = \langle N, h_t h_t^\top \rangle_F = \langle \text{svec}(N), \text{svec}(h_t h_t^\top) \rangle,$$

so this quadratic regression is essentially a linear regression problem in terms of  $\text{svec}(N)$ . A major difficulty in the analysis is to establish persistency of excitation for  $(\text{svec}(h_t h_t^\top))_{t \geq H}$ , meaning that the minimum eigenvalue of the design matrix  $\sum_{t=H}^{T+H-1} \text{svec}(h_t h_t^\top) \text{svec}(h_t h_t^\top)^\top$  grows linearly in the size  $T$  of the data.

A linear lower bound on  $\lambda_{\min}(\sum_{t=H}^{T+H-1} h_t h_t^\top)$  is a known result for the identification of partially observable linear dynamical systems, see the recent overview in (Tsiamis et al., 2022).

In our case, however, elements of  $\text{svec}(h_t h_t^\top)$  are *products* of Gaussians, making the analysis difficult. If  $(h_t)_{t \geq H}$  are independent, which is the case if they are from multiple independent trajectories, the result has been established in (Jadbabaie et al., 2021; Tian et al., 2022). It can also be proved with the matrix Azuma inequality (Tropp, 2012). Here, we need to aggregate dependent data to estimate a set of *stationary* parameters. In sum, the difficulty we face results from both products of Gaussians and the data dependence.

In principle, given enough burn-in time, state  $x_t$ , and hence observation  $y_t$  and truncated history  $h_t$ , converges to the steady-state distributions, and samples with an interval of the order of mixing time are approximately independent (Levin and Peres, 2017). Hence, intuitively, a linear lower bound is viable. However, the bound yielded by such an analysis deteriorates as the system is less stable and the mixing time increases, which is qualitatively incorrect for linear systems. To eschew such dependence, (Simchowitz et al., 2018) introduces the small-ball method. We take the same route, while establishing different arguments to handle the products of Gaussians.

Let us first recall the block martingale small-ball condition.

**Definition 1** (Block martingale small ball (BMSB) condition (Simchowitz et al., 2018, Definition 2.1)). *Let  $(f_t)_{t \geq 1}$  be a stochastic process in  $\mathbb{R}^d$  adapted to filtration  $(\mathcal{F}_t)_{t \geq 1}$ . We say  $(f_t)_{t \geq 1}$  satisfies the  $(k, \Gamma, q)$ -BMSB condition for  $k \in \mathbb{N}^+$ ,  $\Gamma \succcurlyeq 0$  and  $q > 0$ , if for any  $t \geq 1$ , for any fixed unit vector  $v \in \mathbb{R}^d$ ,  $\frac{1}{k} \sum_{i=1}^k \mathbb{P}(|\langle f_{t+i}, v \rangle| \geq \|v\|_\Gamma \mid \mathcal{F}_t) \geq q$  almost surely.*

The key Lemma 1 below shows that  $(\text{svec}(h_t h_t^\top))_{t \geq H}$  satisfies the BMSB condition.

**Lemma 1.** *Let  $h_t = [y_{(t-H+1):t}; u_{(t-H):(t-1)}]$  be the  $H$ -step history at time step  $t$  in system (2.1) with  $u_t \sim \mathcal{N}(0, \sigma_u^2 I)$  for  $t \geq 0$ . Define filtration  $\mathcal{F}_t := \sigma(x_0, v_0, u_0, w_0, v_1, \dots, u_{t-1}, w_{t-1}, v_t)$ . Define  $f_t := \text{svec}(h_t h_t^\top)$ , adapted to  $(\mathcal{F}_t)_{t \geq H}$ . As long as  $H \geq \frac{a \log(\alpha(A^*) \log(T/p))}{\log(\rho(A^*)^{-1})}$ ,  $(f_t)_{t \geq H}$  is  $(k, \gamma^2 I, q)$ -BMSB for  $k = 4H$ ,  $\gamma = \Theta(1/d_h)$ , and  $q = \Theta(1/d_h^2)$ , where  $\Theta(\cdot)$  hides the dependence on dimension-free constants.*

Then, following the analysis in (Simchowitz et al., 2018, Appendix D), we can show that as long as  $T$  is large enough, with high probability,

$$\lambda_{\min} \left( \sum_{t=H}^{T+H-1} f_t f_t^\top \right) = \Omega(\gamma^2 q^2 T) = \Omega(T/d_h^6),$$

which establishes the persistency of excitation.

The full proof of Lemma 1 is deferred to Appendix D.1. Crucial for its proof is the following Lemma 2, which might be of independent interest.

**Lemma 2.** *Let  $x$  be a  $d$ -dimensional zero-mean Gaussian random vector with covariance  $\Sigma$ . Let  $A$  be a  $d \times d$  symmetric matrix with unit Frobenius norm. Then  $\mathbb{E}[|x^\top A x|] \geq a \lambda_{\min}(\Sigma)/d$  for some absolute constant  $a > 0$ .*

## 4.2 Main ideas in proving Theorem 2

Below we sketch the ideas of proving Theorem 2 for CoReDyL.

For the quadratic regression problems (3.2) and (3.6), by establishing persistency of excitation (Lemma 7) and using sub-Weibull martingale concentration (Lemma 6), we can show that

$$\begin{aligned}\|\hat{N} - N^*\|_F &= \mathcal{O}((H(d_y + d_u))^7 d_x^{3/2} T^{-1/2} \log^4(T/p)), \\ \|\hat{N}_1 - N_1^*\|_F &= \mathcal{O}((H(d_y + d_u) + d_x)^7 d_x^{3/2} T^{-1/2} \log^4(T/p));\end{aligned}$$

our choice of  $H$  ensures that the truncation errors are absorbed into these bounds. By the Procrustes-type lemma in (Tu et al., 2016, Lemma 5.4), these two bounds imply bounds of the same order on  $\|\hat{M} - SM^{*'}\|_F$  and  $\|\hat{M}_1 - S_1 M_1^{*'}\|_F$ , respectively, for some orthonormal matrices  $S$  and  $S_1$ . Thus,  $\|\tilde{B} - S_1 B^{*'}\|_2$  is of the same order as the bound on  $\|\hat{N}_1 - N_1^*\|_2$ , and by the perturbation bounds of the Moore-Penrose inverse (Wedin, 1973), so is  $\|\tilde{A} - S_1 A^{*'} S^\top\|_2$ .

**Remark.** Although the next state  $z_{t,1}$  generated by  $(\tilde{A}, \tilde{B})$  yields the correct cumulative cost at the next time step, it mismatches the next state  $\hat{z}_{t+1}$  generated by  $\hat{M}$  using  $h_{t+1}$  by an orthonormal transformation  $S_0 = SS_1^\top$ ; that is,  $z_{t,1}$  approximates  $S_0^\top \hat{z}_{t+1}$  instead of  $\hat{z}_{t+1}$ . This happens because the cost is *invariant* to orthonormal transformations of the latent states, and by only predicting one step into the future, the orthonormal transformations from the two quadratic regressions are not guaranteed to be the same. MuZero bypasses this problem by predicting multiple steps; here we simply calculate this alignment matrix  $S_0$  by Line 5 in Algorithm 2.

As explained, we need to fit the matrix  $S_0 = SS_1^\top$  to ensure the next states generated by  $(\hat{A}, \hat{B})$  align with those by  $\hat{M}$ . With an analysis of *perturbed* linear regression (Lemma 8), we find that  $\|\hat{S}_0 - S_0\|_2 = \mathcal{O}((H(d_y + d_u) + d_x)^{15/2} d_x^{5/2} T^{-1/2} \log^6(T/p))$ . Therefore, the bounds on  $\|\hat{A} - SA^* S^\top\|_2$  and  $\|\hat{B} - SB^*\|_2$  are of the same order. Line 7 in Algorithm 1 requires an analysis of perturbed quadratic regression (Lemma 7), which guarantees that

$$\|\hat{Q} - SQ^* S^\top\|_F = \mathcal{O}((H(d_y + d_u) d_x)^{15/2} T^{-1/2} \log^6(T/p)).$$

Hence,  $\|\hat{A} - SA^* S^\top\|_2$ ,  $\|\hat{B} - SB^*\|_2$  and  $\|\hat{Q} - SQ^* S^\top\|_2$  are all  $\mathcal{O}((H(d_y + d_u) d_x)^{15/2} T^{-1/2} \log^6(d_x T/p))$ .

Lastly, we invoke the result on certainty equivalent LQR in (Mania et al., 2019) to certify the suboptimality gap of  $\hat{K}$  obtained from  $(\hat{A}, \hat{B}, \hat{Q}, R^*)$ . This certainty equivalent controller can also be replaced by a robust one (Dean et al., 2020). Let  $K^{*'}$  denote the optimal controller under the normalized parameterization. Then, for a large enough burn-in period, the certainty-equivalent controller satisfies that  $\|\hat{K} - K^{*'}\|_2$  is of the same order as the system parameter errors. Note that policy  $K^* M^*$  in the original system is independent of the latent model parameterization, and we have that  $\|\hat{K} \hat{M} - K^* M^*\|_2$  is of the same order.

## 5 Numerical results

Although this work is theoretical in nature, we conduct preliminary numerical experiments by testing CoREL on a mass-spring-damper system, with mass 1 kg, stiffness 10 N/m and damping coefficient 100 kg/s. We discretize time with the forward Euler method using 0.01 s intervals. The state of the system is the position and the velocity of the mass, with the observation being the position only. The scalar control input is the external force. We take

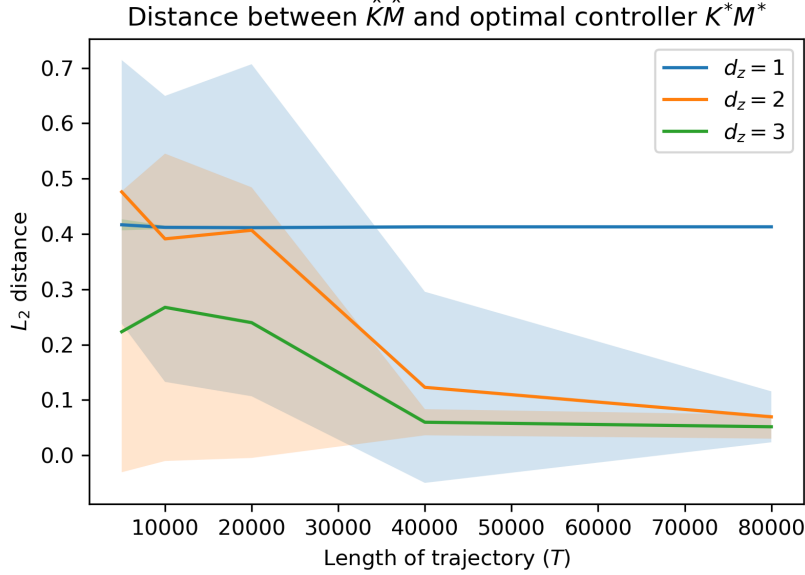


Figure 1: Optimality of the learned policy versus the length of the trajectory. Each dot averages the results of five independent runs.

$Q = I$  and  $R = 1$  in the cost function. We set  $H = 5$  and measure the performance by the distance between  $\hat{K}\hat{M}$  and  $K^*M^*$ . In all experiments, we warm up the system by 1000 steps.

The simulation results are summarized in Figure 1. In general, we observe the decrease of the distance to the optimal policy as the length of the trajectory increases. The system has  $d_x = 2$ . Besides  $d_z = 2$ , we experiment with  $d_z = 1$  and  $d_z = 3$  to examine the impact of dimension mismatch. In this example, we see that  $d_z = 1$  does not produce a meaningful policy and  $d_z = 3$  performs better than the  $d_z = 2$  baseline. We note that these results are preliminary, and this “blessing of overparameterization” in the latent state space is worth further investigation.

## 6 Conclusion and future work

We studied direct latent model learning for solving unknown infinite-horizon time-invariant LQG control. We established finite-sample guarantees for two methods, which differ in whether the latent dynamics is learned explicitly or implicitly, with the latter being closer to that used in MuZero (Schrittwieser et al., 2020). For MuZero-style latent model learning, our analysis identifies a coordinate misalignment problem in the latent state space, suggesting the value of *multi-step* future prediction. A limitation of this work is that we only consider state representation based on truncated histories, i.e., frame stacking, as used in MuZero; the *recursive form* of the representation function, as in the Kalman filter, is also used empirically (Ha and Schmidhuber, 2018; Hafner et al., 2019a), and might be worth further investigation.

Many questions remain to be answered in representation learning for control. Provable



generalization of direct latent model learning to nonlinear observation channels or dynamics is a natural consideration. Moreover, with the ubiquity of visual perception in real-world control systems, what if we have a time-varying observation function or multiple observation functions, modeling images taken from different angles? In reality, most of the time we do not have a well-defined cost function; learning task-relevant state representations from demonstrations is another intriguing direction.

## Acknowledgements

This work was supported in part by the NSF TRIPODS program (award number DMS-2022448).

## References

- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. *Advances in neural information processing systems*, 24, 2011.
- Dimitri Bertsekas. *Dynamic Programming and Optimal Control: Volume I*, volume 1. Athena Scientific, 2012.
- Sarah Dean, Horia Mania, Nikolai Matni, Benjamin Recht, and Stephen Tu. On the sample complexity of the linear quadratic regulator. *Foundations of Computational Mathematics*, 20(4): 633–679, 2020.
- Xiequan Fan, Ion Grama, and Quansheng Liu. Large deviation exponential inequalities for supermartingales. 2012.
- Xiequan Fan, Ion Grama, and Quansheng Liu. Deviation inequalities for martingales with applications. *Journal of Mathematical Analysis and Applications*, 448(1):538–566, 2017.
- Xiang Fu, Ge Yang, Pulkit Agrawal, and Tommi Jaakkola. Learning task informed abstractions. In *International Conference on Machine Learning*, pages 3480–3491. PMLR, 2021.
- David Ha and Jürgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019a.
- Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James Davidson. Learning latent dynamics for planning from pixels. In *International conference on machine learning*, pages 2555–2565. PMLR, 2019b.
- Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering Atari with discrete world models. *arXiv preprint arXiv:2010.02193*, 2020.
- Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models. *arXiv preprint arXiv:2301.04104*, 2023.

- Ali Jadbabaie, Horia Mania, Devavrat Shah, and Suvrit Sra. Time varying regression with hidden linear dynamics. *arXiv preprint arXiv:2112.14862*, 2021.
- N Komaroff. Iterative matrix bounds and computational solutions to the discrete algebraic Riccati equation. *IEEE Transactions on Automatic Control*, 39(8):1676–1678, 1994.
- Sahin Lale, Kamyar Azizzadenesheli, Babak Hassibi, and Anima Anandkumar. Logarithmic regret bound in partially observable linear dynamical systems. *Advances in Neural Information Processing Systems*, 33:20876–20888, 2020.
- Sahin Lale, Kamyar Azizzadenesheli, Babak Hassibi, and Anima Anandkumar. Adaptive control and regret minimization in linear quadratic Gaussian (LQG) setting. In *2021 American Control Conference (ACC)*, pages 2517–2522. IEEE, 2021.
- Alex Lamb, Riashat Islam, Yonathan Efroni, Aniket Didolkar, Dipendra Misra, Dylan Foster, Lekan Molu, Rajan Chari, Akshay Krishnamurthy, and John Langford. Guaranteed discovery of controllable latent states with multi-step inverse models. *arXiv preprint arXiv:2207.08229*, 2022.
- David A Levin and Yuval Peres. *Markov chains and mixing times*, volume 107. American Mathematical Soc., 2017.
- Horia Mania, Stephen Tu, and Benjamin Recht. Certainty equivalence is efficient for linear quadratic control. *Advances in Neural Information Processing Systems*, 32, 2019.
- Nikolai Matni and Stephen Tu. A tutorial on concentration bounds for system identification. In *2019 IEEE 58th Conference on Decision and Control (CDC)*, pages 3741–3749. IEEE, 2019.
- Zakaria Mhammedi, Dylan J Foster, Max Simchowitz, Dipendra Misra, Wen Sun, Akshay Krishnamurthy, Alexander Rakhlin, and John Langford. Learning the linear quadratic regulator from nonlinear observations. *Advances in Neural Information Processing Systems*, 33:14532–14543, 2020.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing Atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- Junhyuk Oh, Satinder Singh, and Honglak Lee. Value prediction network. *Advances in neural information processing systems*, 30, 2017.
- 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.
- Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration by self-supervised prediction. In *International conference on machine learning*, pages 2778–2787. PMLR, 2017.

- Kathrin Schacke. On the Kronecker product. *Master's Thesis, University of Waterloo*, 2004.
- Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering Atari, Go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.
- 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.
- David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of Go without human knowledge. *nature*, 550(7676):354–359, 2017.
- David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dharmashan Kumaran, Thore Graepel, et al. A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play. *Science*, 362(6419):1140–1144, 2018.
- 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.
- Max Simchowitz, Ross Boczar, and Benjamin Recht. Learning linear dynamical systems with semi-parametric least squares. In *Conference on Learning Theory*, pages 2714–2802. PMLR, 2019.
- Max Simchowitz, Karan Singh, and Elad Hazan. Improper learning for non-stochastic control. In *Conference on Learning Theory*, pages 3320–3436. PMLR, 2020.
- Jayakumar Subramanian, Amit Sinha, Raihan Seraj, and Aditya Mahajan. Approximate information state for approximate planning and reinforcement learning in partially observed systems. *arXiv preprint arXiv:2010.08843*, 2020.
- Yi Tian, Kaiqing Zhang, Russ Tedrake, and Suvrit Sra. Can direct latent model learning solve linear quadratic Gaussian control? *arXiv preprint arXiv:2212.14511*, 2022.
- Joel A Tropp. User-friendly tail bounds for sums of random matrices. *Foundations of computational mathematics*, 12:389–434, 2012.
- Anastasios Tsiamis, Ingvar Ziemann, Nikolai Matni, and George J Pappas. Statistical learning theory for control: A finite sample perspective. *arXiv preprint arXiv:2209.05423*, 2022.
- Stephen Tu, Ross Boczar, Max Simchowitz, Mahdi Soltanolkotabi, and Ben Recht. Low-rank solutions of linear matrix equations via Procrustes flow. In *International Conference on Machine Learning*, pages 964–973. PMLR, 2016.
- Per-Åke Wedin. Perturbation theory for pseudo-inverses. *BIT Numerical Mathematics*, 13:217–232, 1973.

- Weirui Ye, Shaohuai Liu, Thanard Kurutach, Pieter Abbeel, and Yang Gao. Mastering Atari games with limited data. *Advances in Neural Information Processing Systems*, 34:25476–25488, 2021.
- Amy Zhang, Rowan McAllister, Roberto Calandra, Yarin Gal, and Sergey Levine. Learning invariant representations for reinforcement learning without reconstruction. *arXiv preprint arXiv:2006.10742*, 2020.
- Huiming Zhang and Haoyu Wei. Sharper sub-Weibull concentrations. *Mathematics*, 10(13):2252, 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.
- Yang Zheng, Luca Furieri, Maryam Kamgarpour, and Na Li. Sample complexity of linear quadratic Gaussian (LQG) control for output feedback systems. In *Learning for Dynamics and Control*, pages 559–570. PMLR, 2021.

## A Additional related work

(Oymak and Ozay, 2019) studies the identification of partially observable linear dynamical systems from a single trajectory, which presents a finite-sample analysis of identifying the Markov parameter and a perturbation analysis of the Ho-Kalman algorithm. (Simchowitz et al., 2019) relaxes the stability requirement to marginal stability by using prefiltered least squares to identify the Markov parameter. The method in (Zheng and Li, 2020) applies to unstable systems but requires multiple trajectories. Since the Markov parameter maps control input histories to observations, these methods do not work with costs and use the Markov parameter as an intermediate step to identify the system. By contrast, our methods, entirely driven by the costs and closely connected with empirical methods, directly learn the representation function and the latent model. Directly learning the latent model connects our work to the identification of fully observable linear dynamical systems. (Simchowitz et al., 2018) introduces small-ball conditions to handle dependent data and characterizes the statistical rates for stable and unstable systems, both proving to be useful for our analysis.

Online control of partially observable linear dynamical systems is considered in (Lale et al., 2020, 2021) for stochastic noises and in (Simchowitz et al., 2020) for nonstochastic noises. (Zheng et al., 2021) considers end-to-end sample complexity and is closest to our setup. All these methods rely on the estimation of Markov parameters. For a discussion of the literature in more details and breadth, we refer the reader to the recent survey (Tsiamis et al., 2022).

## B Additional discussions on CoSysId

In CoSysId (Algorithm 2), the covariates of quadratic regression in (3.6) are  $([h_t; u_t])_{t \geq H}$ . One may wonder if we can pursue an alternative approach by fixing  $M$  to be  $\hat{M}$ , and using  $([\hat{z}_t; u_t])_{t \geq H}$  as covariates, which have a much lower dimension, though the two quadratic regressions cannot be solved in parallel anymore. Specifically, the new quadratic regression we need to solve is given by

$$\hat{N}_2, \hat{b}_2 \in \underset{N_2=N_2^\top, b_2}{\operatorname{argmin}} \sum_{t=H}^{T+H-1} (\|[\hat{z}_t; u_t]\|_{N_2}^2 + b_2 - \bar{c}_{t+1})^2,$$

where  $\hat{z}_t = \hat{M}h_t$  is an approximation of  $Sz_t^*$ . The ground truth for  $\hat{N}_2$  is  $N_2^* = [SA^*S^\top, SB^*]^\top [SA^*S^\top, SB^*]$ , so its approximate factorization recovers  $[S_3A^*S^\top, S_3B^*]$  for some orthonormal matrix  $S_3$ . In a similar way to CoSysId, we still need to fit an alignment matrix  $S_2 = SS_3^\top$  to align the coordinates. Let  $\tilde{A}, \tilde{B}$  denote the system parameters recovered from  $\hat{N}_2$ . The linear regression we now need to solve is from  $([\tilde{A}, \tilde{B}][\hat{z}_t; u_t])_{t=H}^{T+H-1}$  to  $(\hat{z}_{t+1})_{t=H}^{T+H-1}$ . However, without further assumptions,  $[A^*, B^*]$  does not necessarily have full row-rank, and hence, neither does  $[\tilde{A}, \tilde{B}]$ , in which case recovering the entire  $S_2$  is impossible.

## C Auxiliary results

### C.1 Basic inequalities involving random vectors

**Lemma 3** ((Tian et al., 2022, Lemma 3)). *Let  $x$  and  $y$  be random vectors defined on the same probability space. Then,  $\|\text{Cov}(x, y)\|_2 \leq \|\text{Cov}(x)^{1/2}\|_2 \|\text{Cov}(y)^{1/2}\|_2$ .*

**Lemma 4** ((Tian et al., 2022, Lemma 2)). *Let  $x, y$  be random vectors of dimensions  $d_x, d_y$ , respectively, defined on the same probability space. Then,  $\|\text{Cov}([x; y])\|_2 \leq \|\text{Cov}(x)\|_2 + \|\text{Cov}(y)\|_2$ .*

### C.2 Lower bound about Gaussian quadratic forms

**Lemma 5.** *Let  $z_1, z_2, \dots, z_d$  be independent standard Gaussian random variables. Let  $v = [v_1, v_2, \dots, v_d]^\top \in \mathbb{R}^d$  be a unit vector. There exists an absolute constant  $a > 0$ , such that for any such  $v$ ,  $\mathbb{E}[\|\sum_{i=1}^d v_i z_i^2\|] \geq a/d$ .*

*Proof of Lemma 5.* Let  $\text{sign}(\cdot)$  denote the sign function. Let  $\mathcal{I}^+ := \{i : \text{sign}(v_i) = 1, 1 \leq i \leq d\}$  and  $\mathcal{I}^- := \{i : \text{sign}(v_i) = -1, 1 \leq i \leq d\}$  be the index sets of positive and negative values. Then,

$$\mathbb{E}\left[\left\|\sum_{i=1}^d v_i z_i^2\right\|\right] = \mathbb{E}\left[\left\|\sum_{i=1}^d |v_i| \text{sign}(v_i) z_i^2\right\|\right] = \mathbb{E}\left[\left\|\sum_{i \in \mathcal{I}^+} |v_i| z_i^2 - \sum_{j \in \mathcal{I}^-} |v_j| z_j^2\right\|\right].$$

For a given  $v$ , since  $(z_i^2)_{i=1}^d$  have identical distributions,  $\mathbb{E}\left[\left\|\sum_{i \in \mathcal{I}^+} |v_i| z_i^2 - \sum_{j \in \mathcal{I}^-} |v_j| z_j^2\right\|\right]$  has the same value under permutations of  $(v_i)_{i \in \mathcal{I}^+}$  and  $(v_j)_{j \in \mathcal{I}^-}$ . Summing over all the permutations of  $(v_i)_{i \in \mathcal{I}^+}$  and  $(v_j)_{j \in \mathcal{I}^-}$  gives

$$d \mathbb{E}\left[\left\|\sum_{i=1}^d v_i z_i^2\right\|\right] \geq \mathbb{E}\left[\left|\left(\sum_{i \in \mathcal{I}^+} |v_i|\right) \sum_{i \in \mathcal{I}^+} z_i^2 - \left(\sum_{j \in \mathcal{I}^-} |v_j|\right) \sum_{j \in \mathcal{I}^-} z_j^2\right|\right].$$

Hence,

$$\mathbb{E}\left[\left\|\sum_{i=1}^d v_i z_i^2\right\|\right] \geq \frac{1}{d} \left(\sum_{i=1}^d |v_i|\right) \mathbb{E}\left[\left\|\sum_{i=1}^d \text{sign}(v_i) z_i^2\right\|\right]$$

Since  $\sum_{i=1}^d |v_i| \geq (\sum_{i=1}^d v_i^2)^{1/2} = 1$ , we have

$$\mathbb{E}\left[\left\|\sum_{i=1}^d v_i z_i^2\right\|\right] \geq \frac{1}{d} \inf_{w \in \{\pm 1\}^d} \mathbb{E}\left[\left\|\sum_{i=1}^d w_i z_i^2\right\|\right].$$

It remains to lower bound  $\inf_{w \in \{\pm 1\}^d} \mathbb{E}\left[\left\|\sum_{i=1}^d w_i z_i^2\right\|\right]$ . Let  $p$  denote the number of  $+1$ 's and  $q$  denote the number of  $-1$ 's in  $w$ , such that  $p + q = n$ . If  $p \neq q$ , by Jensen's inequality,

$$\mathbb{E}\left[\left\|\sum_{i=1}^d w_i z_i^2\right\|\right] \geq \mathbb{E}\left[\sum_{i=1}^{|p-q|} z_i^2\right] \geq \mathbb{E}[z_1^2] = \Omega(1).$$

If  $p = q$ , again, an application of Jensen's inequality yields

$$\mathbb{E}\left[\left\|\sum_{i=1}^d w_i z_i^2\right\|\right] \geq \mathbb{E}[|z_1^2 - z_2^2|] = \Omega(1).$$

Overall, we have  $\inf_{w \in \{\pm 1\}^d} \mathbb{E}\left[\left\|\sum_{i=1}^d w_i z_i^2\right\|\right] = \Omega(1)$ . Hence,  $\mathbb{E}\left[\left\|\sum_{i=1}^d v_i z_i^2\right\|\right] \geq \Omega(1/d)$ .  $\square$

Below we use Lemma 5 to prove Lemma 2.

*Proof of Lemma 2.* Let  $y := \Sigma^{-1/2}x$ . Then  $y$  is a standard Gaussian random vector, and  $x^\top Ax = y^\top \Sigma^{1/2} A \Sigma^{1/2} y$ . Let  $U^\top \Lambda U$  be the eigenvalue decomposition of  $\Sigma^{1/2} A \Sigma^{1/2}$ . Then,

$$\mathbb{E}[|x^\top Ax|] = \mathbb{E}[|y^\top U^\top \Lambda U y|] = \mathbb{E}[|z^\top \Lambda z|],$$

where  $z := Uy$  is still a standard Gaussian random vector. Since

$$\|\Lambda\|_F \stackrel{(i)}{=} \|U^\top \Lambda U\|_F = \|\Sigma^{1/2} A \Sigma^{1/2}\|_F \geq \lambda_{\min}(\Sigma) \|A\|_F = \lambda_{\min}(\Sigma),$$

where (i) is due to the unitary invariance of the Frobenius norm, we have

$$\inf_{\|A\|_F=1} \mathbb{E}[|x^\top Ax|] \geq \inf_{\|\Lambda\|_F \geq \lambda_{\min}(\Sigma)} \mathbb{E}[|z^\top \Lambda z|] \stackrel{(i)}{\geq} a \lambda_{\min}(\Sigma) / d,$$

where (i) is due to Lemma 5.  $\square$

### C.3 Sum of sub-Weibull martingale difference sequences

To upper bound  $\sum_{t=H}^{T+H-1} f_t e_t$  in the analysis of quadratic regression (see §D), one possible approach (Tsiamis et al., 2022) is to use bounds for self-normalized martingales (Abbasi-Yadkori et al., 2011), but the standard self-normalized martingale lemma assumes the noises  $(e_t)_{t \geq H}$  to be sub-Gaussian. (Fan et al., 2012, 2017) study the sum of martingale difference sequences with sub-Weibull distributions, based on which we prove Lemma 6.

**Lemma 6.** Let  $(\eta_t)_{t \geq 1}$  be a martingale difference sequence adapted to filtration  $(\mathcal{F}_t)_{t \geq 1}$ . Assume  $\eta_t \mid \mathcal{F}_{t-1}$  is  $\theta$ -sub-Weibull with  $\|\eta_t \mid \mathcal{F}_{t-1}\|_{\psi_\theta} \leq K$ . Then with probability at least  $1 - p$ , there exist absolute constants  $c, c' > 0$ , such that as long as  $n \geq c$ ,

$$\sum_{t=1}^T \eta_t \leq c' K \sqrt{T} (\log(T/p))^{1+\theta^{-1}}.$$

*Proof.* By the definition of sub-Weibull distributions,  $\mathbb{E}[\exp(|\eta_t|/K)^\theta \mid \mathcal{F}_{t-1}] \leq 2$ . Define  $\epsilon_t = \eta_t/K$ . Then by the properties of sub-Weibull distributions,  $\mathbb{E}[\epsilon_t^2 \mid \mathcal{F}_{t-1}] \leq a$ , for some absolute constant  $a > 0$ . Hence,  $(\epsilon_t)_{t \geq 1}$  satisfies the assumptions required in (Fan et al., 2017, Theorem 3.2) for  $\alpha = \theta/(\theta + 1)$ . Taking  $(\phi_t)_{t \geq 1}$  in (Fan et al., 2017, Theorem 3.2) to be ones, we have

$$\mathbb{P}\left(\sum_{t=1}^T \epsilon_t \geq x\sqrt{T}\right) \leq \exp\left(-\frac{x^2}{2(c + x^{1+\frac{1}{\theta+1}}/3)}\right) + 2T \exp\left(-x^{\frac{\theta}{\theta+1}}\right).$$

Note that

$$\frac{x^2}{2(c + x^{1+\frac{1}{\theta+1}}/3)} = \frac{x^{\frac{\theta}{\theta+1}}}{2/3 + 2cx^{-1-\frac{1}{\theta+1}}} \geq x^{\frac{\theta}{\theta+1}},$$

if  $2cx^{-1-\frac{1}{\theta+1}} \leq 1/3$ , that is,  $x \geq (6c)^{\frac{\theta+1}{\theta+2}}$ . Then, as long as  $x \geq 6c$ ,

$$\mathbb{P}\left(\sum_{t=1}^T \epsilon_t \geq x\sqrt{T}\right) \leq 3T \exp\left(-x^{\frac{\theta}{\theta+1}}\right).$$



Hence,

$$\mathbb{P}\left(\sum_{t=1}^T \epsilon_t \geq 6c + \sqrt{T}(\log(3T/p))^{1+\theta^{-1}}\right) \leq p.$$

Therefore, there exists an absolute constant  $c' > 0$ , such that if  $T \geq \max(36c^2, 3)$ , then with probability at least  $1 - p$ ,

$$\sum_{t=1}^T \epsilon_t \leq c'K\sqrt{T}(\log(T/p))^{1+\theta^{-1}}.$$

□

#### C.4 Proposition on multi-step cumulative costs

The following proposition, taken from (Tian et al., 2022, Proposition 3) and adapted to our LTI setting, is important for analyzing CoREL and CoREDyL.

**Proposition 3.** Define filtration  $\mathcal{F}_t := \sigma(x_0, v_0, u_0, w_0, v_1, \dots, u_{t-1}, w_{t-1}, v_t)$ . Let  $z_t^{*'} = \hat{x}_t'$  be the state estimates generated by the Kalman filter under the normalized parameterization, adapted to  $(\mathcal{F}_t)_{t \geq 0}$ . If we apply  $u_t \sim \mathcal{N}(0, \sigma_u^2 I)$ , then for any  $t \geq 0$ ,

$$\bar{c}_t := \sum_{\tau=t}^{t+d_x-1} (c_\tau - \|u_\tau\|_{R^*}^2) = \|z_t^{*'}\|^2 + b' + e'_t,$$

where  $b' = \mathcal{O}(d_x)$  is a positive constant, and  $e'_t \mid \mathcal{F}_{t-1}$  is a zero-mean subexponential random variable with  $\|e'_t \mid \mathcal{F}_{t-1}\|_{\psi_1} = \mathcal{O}(d_x^{3/2})$ .

## D Quadratic regression bound

The following quadratic regression bound is at the core of proving Theorem 2. Its proof builds on a new persistency of excitation result (Lemma 1) and the concentration of sub-Weibull martingale different sequences (Lemma 6).

**Lemma 7.** Let  $(h_t^*)_{t \geq 1}$  be a sequence of  $d$ -dimensional Gaussian random vectors adapting to filtration  $(\mathcal{F}_t)_{t \geq 1}$  with  $\|\mathbb{E}[h_t^* (h_t^*)^\top]\|_2^{1/2} \leq \sigma$ . Define random variable  $c_t = (h_t^*)^\top N h_t^* + b^* + e_t$ , where  $e_t \mid \mathcal{F}_{t-1}$  is zero-mean, subexponential with  $\|e_t \mid \mathcal{F}_{t-1}\|_{\psi_1} \leq E$ . Define  $h_t = h_t^* + \delta_t$ , where the Gaussian noise vector  $\delta_t$  can be correlated with  $h_t^*$  and satisfy  $\|\mathbb{E}[\delta_t \delta_t^\top]\|_2^{1/2} \leq \epsilon \leq \sigma$ . Assume that  $(\text{vec}(h_t^* h_t^{*\top}))_{t \geq 1}$  satisfies the  $(k, \gamma^2 I, q)$ -BMSB condition and  $\epsilon \leq a_0 \gamma \sigma^{-1} d^{-1} \log^{-2}(T/p)$  for some absolute constant  $a_0 > 0$ . Consider

$$(\hat{N}, \hat{b}) \in \underset{N=N^\top, b}{\operatorname{argmin}} \sum_{t=1}^T (c_t - \|h_t\|_N^2 - b)^2. \quad (\text{D.1})$$

Then, as long as  $T \geq a_1 k d^5 \log(d/p)$  for some dimension-free constant  $a_1 > 0$ , we have that with probability at least  $1 - p$ ,

$$\|\hat{N} - N^*\|_F = \mathcal{O}(\epsilon(\gamma q)^{-1} d \log^2(T/p) + (\gamma q)^{-2} d E T^{-1/2} \log^4(T/p)),$$

where  $\sigma$ ,  $E$  and  $\|N^*\|_2$  are problem-dependent constants hidden in  $\mathcal{O}(\cdot)$ .

*Proof.* Let  $f_t := \text{svec}(h_t h_t^\top)$  and  $F := [f_1, f_2, \dots, f_T]^\top$  be the  $T \times \frac{d(d+1)}{2}$  matrix whose  $t$ th row is  $f_t^\top$ . Define  $f_t^*$  and  $F^*$  similarly. Using  $f_t$ , we can rewrite the regression (D.1) as

$$(\text{svec}(\hat{N}), \hat{b}) \in \underset{\text{svec}(N), b}{\text{argmin}} \sum_{t=1}^T (c_t - f_t^\top \text{svec}(N) - b)^2. \quad (\text{D.2})$$

Solving regression (D.2) for  $N$  and  $b$ , and substituting in the expression of  $c_t$ , we obtain

$$F^\top F \text{svec}(\hat{N}) = F^\top F^* \text{svec}(N^*) + (b^* - \hat{b}) F^\top 1_T + F^\top \varepsilon, \quad (\text{D.3})$$

$$\hat{b} = b^* + \frac{1}{T} (1_T^\top F^* \text{svec}(N^*) - 1_T^\top F \text{svec}(\hat{N}) + 1_T^\top \varepsilon), \quad (\text{D.4})$$

where  $\varepsilon$  denotes the vector whose  $t$ -th element is  $e_t$ . By substituting (D.4) into (D.3), we have

$$F^\top \left( I_T - \frac{1_T 1_T^\top}{T} \right) F \text{svec}(\hat{N}) = F^\top \left( I_T - \frac{1_T 1_T^\top}{T} \right) F^* \text{svec}(N^*) + F^\top \left( I_T - \frac{1_T 1_T^\top}{T} \right) \varepsilon.$$

To deal with  $J_T = (I_T - 1_T 1_T^\top / T)$ , we adopt the same reparameterization trick in (Tian et al., 2022), by defining  $\tilde{F} := J_T^{1/2} F$  and  $\tilde{F}^* := J_T^{1/2} F^*$ . Since  $J_T$  is a positive definite matrix with  $T - 1$  eigenvalues being one and the other being  $1 - 1/T$ , the eigenvalues of  $\tilde{F}^\top \tilde{F}$  do not differ much from those of  $F^\top F$ . In particular,  $\lambda_{\min}(\tilde{F}^\top \tilde{F}) \geq (1 - 1/T) \lambda_{\min}(F^\top F)$ . With this reparameterization and subtracting  $\tilde{F}^\top \tilde{F} \text{svec}(N^*)$  from both sides, we have

$$\tilde{F}^\top \tilde{F} \text{svec}(\hat{N} - N^*) = \tilde{F}^\top (\tilde{F}^* - \tilde{F}) \text{svec}(N^*) + \tilde{F}^\top J_T^{1/2} \varepsilon. \quad (\text{D.5})$$

Below we bound  $\|(\tilde{F}^* - \tilde{F}) \text{svec}(N^*)\|$ ,  $\|\tilde{F}^\top J_T^{1/2} \varepsilon\|$ , and  $\lambda_{\min}(\tilde{F}^\top \tilde{F})$  separately.

To bound  $\|(\tilde{F}^* - \tilde{F}) \text{svec}(N^*)\|$  we apply similar techniques in (Mhammedi et al., 2020; Tian et al., 2022), by first noting that

$$\begin{aligned} \|(F^* - F) \text{svec}(N^*)\|^2 &= \sum_{t=1}^T \left\langle \text{svec}(h_t^* (h_t^*)^\top) - \text{svec}(h_t h_t^\top), \text{svec}(N^*) \right\rangle \\ &= \sum_{t=1}^T \left\langle h_t^* (h_t^*)^\top - h_t h_t^\top, N^* \right\rangle_F \\ &= \|N^*\|_2^2 \sum_{t=1}^T \|h_t^* (h_t^*)^\top - h_t h_t^\top\|_*^2 \\ &\stackrel{(i)}{\leq} 2 \|N^*\|_2^2 \sum_{t=1}^T \|h_t^* (h_t^*)^\top - h_t h_t^\top\|_F^2 = 2 \|N^*\|_2^2 \|F^* - F\|_F^2, \end{aligned}$$

where (i) follows from the fact that the matrix  $h_t^* (h_t^*)^\top - h_t h_t^\top$  has at most rank two.

$$\begin{aligned} \|F^* - F\|_F^2 &= \sum_{t=1}^T \|h_t^* (h_t^*)^\top - h_t h_t^\top\|_F^2 \\ &= \sum_{t=1}^T \|h_t^* (h_t^* - h_t)^\top + (h_t^* - h_t) h_t^\top\|_F^2 \\ &= \sum_{t=1}^T \|h_t^* \delta_t^\top + \delta_t h_t^\top\|_F^2 \leq \sum_{t=1}^T 2(\|h_t^*\|^2 + \|h_t\|^2) \|\delta_t\|^2. \end{aligned}$$

With probability at least  $1 - p$ , for all  $1 \leq t \leq T$ ,

$$\|h_t^*\| = \mathcal{O}(\sigma d^{1/2} \log(T/p)), \quad \|h_t\| = \mathcal{O}(\sigma d^{1/2} \log(T/p)), \quad \|\delta_t\| = \mathcal{O}(\epsilon d^{1/2} \log(T/p)).$$

Hence,

$$\|F^* - F\|_F^2 = \mathcal{O}(\epsilon^2 \sigma^2 d^2 T \log^4(T/p)).$$

It thus follows that

$$\|(\tilde{F}^* - \tilde{F})\text{svec}(N^*)\| = \|J_T^{1/2}(F^* - F)\text{svec}(N^*)\| = \mathcal{O}(\epsilon \sigma d \|N^*\|_2 T^{1/2} \log^2(T/p)).$$

Next, we handle  $\|\tilde{F}^\top J_T^{1/2} \varepsilon\|$  by

$$\|\tilde{F}^\top J_T^{1/2} \varepsilon\| = \left\| F^\top \left( I - \frac{1_T 1_T^\top}{T} \right) \varepsilon \right\| \leq \|F^\top \varepsilon\| + \frac{1}{T} \|F^\top 1_T 1_T^\top \varepsilon\| = \left\| \sum_{t=1}^T f_t e_t \right\| + \frac{1}{T} \left\| \left( \sum_{t=1}^T f_t \right) \left( \sum_{t=1}^T e_t \right) \right\|.$$

Since  $\|[f_t]_i \mid \mathcal{F}_{t-1}\|_{\psi_1} = \mathcal{O}(\sigma^2)$  and  $\|e_t \mid \mathcal{F}_{t-1}\|_{\psi_1} \leq E$ , the product  $[f_t]_i e_t \mid \mathcal{F}_{t-1}$  is  $\frac{1}{2}$ -sub-Weibull, with the sub-Weibull norm being  $\mathcal{O}(\sigma^2 E)$ . By Lemma 6, with probability at least  $1 - p$ ,

$$\sum_{t=1}^T [f_t]_i e_t = \mathcal{O}(\sigma^2 E T^{1/2} \log^3(T/p)).$$

Then, since  $\sum_{t=H}^{T+H-1} f_t e_t$  has  $d(d+1)/2$  components,

$$\left\| \sum_{t=H}^{T+H-1} f_t e_t \right\| = \mathcal{O}(\sigma^2 d E T^{1/2} \log^3(T/p)).$$

On the other hand, with probability at least  $1 - p$ ,

$$\left\| \sum_{t=1}^T f_t \right\| \leq \sum_{t=1}^T \|h_t h_t^\top\|_F = \sum_{t=1}^T \|h_t\|^2 = \mathcal{O}(\sigma^2 d T \log^2(T/p)),$$

and by Lemma 6, with probability at least  $1 - p$ ,

$$\left\| \sum_{t=1}^T e_t \right\| = \mathcal{O}(E T^{1/2} \log^2(T/p)).$$

Hence,

$$\|\tilde{F}^\top J_T^{1/2} \varepsilon\| = \mathcal{O}(\sigma^2 d E T^{1/2} \log^3(T/p) + \sigma^2 d \log^2(T/p) E T^{1/2} \log^2(T/p)) = \mathcal{O}(\sigma^2 d E T^{1/2} \log^4(T/p)).$$

The remaining term to deal with is  $\lambda_{\min}(\tilde{F}^\top \tilde{F}) = \Omega(F^\top F)$ , which we achieve by showing  $(f_t)_{t \geq 1}$  is BMSB. By our assumption,  $(f_t^*)_{t \geq 1}$  is  $(k, \gamma^2 I, q)$ -BMSB, meaning that for any fixed unit vector  $v \in \mathbb{R}^{\frac{d(d+1)}{2}}$ , it holds almost surely that

$$\frac{1}{k} \sum_{i=1}^k \mathbb{P}(|\langle f_{t+i}^*, v \rangle| \geq \gamma \mid \mathcal{F}_t) \geq q.$$

For any fixed unit vector  $v \in \mathbb{R}^{\frac{d(d+1)}{2}}$ , we have

$$|\langle f_t, v \rangle| = |\langle f_t^*, v \rangle + \langle f_t - f_t^*, v \rangle| \geq |\langle f_t^*, v \rangle| - |\langle f_t - f_t^*, v \rangle| \geq |\langle f_t^*, v \rangle| - \|f_t - f_t^*\|.$$

For all  $1 \leq t \leq T$ , since  $\|f_t - f_t^*\| = \|h_t h_t^\top - h_t^* (h_t^*)^\top\|_F = \mathcal{O}(\epsilon \sigma d \log^2(T/p))$ , there exists an absolute constant  $a_0 > 0$ , such that as long as  $\epsilon \leq \frac{a_0 \gamma}{\sigma d \log^2(T/p)}$ ,  $\|f_t - f_t^*\| \leq \gamma/2$ . It follows that

$$\frac{1}{k} \sum_{i=1}^k \mathbb{P}(|\langle f_{t+i}, v \rangle| \geq \gamma/2 \mid \mathcal{F}_t) \geq \frac{1}{k} \sum_{i=1}^k \mathbb{P}(|\langle f_{t+i}^*, v \rangle| \geq \gamma \mid \mathcal{F}_t) \geq q,$$

which means that  $(f_t)_{1 \leq t \leq T}$  is  $(k, \gamma^2 I/4, q)$ -BMSB. Following the analysis in (Simchowitz et al., 2018, Appendix D), we can show that for a given  $p \in (0, 1)$ , as long as  $T \geq a_1 k d^5 \log(d/p)$  for some dimension-free constant  $a_1 > 0$ , then with probability at least  $1 - p$ , we have

$$\lambda_{\min}\left(\sum_{t=1}^T f_t f_t^\top\right) = \Omega(\gamma^2 q^2 T).$$

Hence, we have  $\lambda_{\min}(\tilde{F}^\top \tilde{F}) = \Omega(\gamma^2 q^2 T)$ .

Finally, recall that by (D.5) we have

$$\text{svec}(\hat{N} - N^*) = \tilde{F}^\top (\tilde{F}^* - \tilde{F}) \text{svec}(N^*) + (\tilde{F}^\top \tilde{F})^{-1} \tilde{F}^\top J_T^{1/2} \varepsilon,$$

by combining the individual bounds proved above, we show that there exists some absolute constant  $a > 0$ , such that as long as  $T \geq a_1 k d^5 \log(d/p)$ , with probability at least  $1 - p$ ,

$$\begin{aligned} \|\hat{N} - N^*\|_F &= \|\text{svec}(\hat{N} - N^*)\|_2 \leq \lambda_{\min}(\tilde{F}^\top \tilde{F})^{-1/2} \|(\tilde{F}^* - \tilde{F}) \text{svec}(N^*)\| + \lambda_{\min}(\tilde{F}^\top \tilde{F})^{-1} \|\tilde{F}^\top J_T^{1/2} \varepsilon\| \\ &= \mathcal{O}(\epsilon(\gamma q)^{-1} d \log^2(T/p) + (\gamma q)^{-2} d E T^{-1/2} \log^4(T/p)), \end{aligned}$$

where  $\sigma$  and  $\|N^*\|_2$  are problem-dependent constants hidden in  $\mathcal{O}(\cdot)$ .

As a side remark, from (D.4), we obtain

$$|\hat{b} - b^*| = \left| \frac{1}{T} \mathbf{1}_T^\top (F^* - F) \text{svec}(N^*) \right| + \left| \frac{1}{T} \mathbf{1}_T^\top F \text{svec}(N^* - \hat{N}) \right| + \left| \frac{1}{T} \mathbf{1}_T^\top \varepsilon \right|.$$

Substituting in the individual bounds we have derived and using  $\|\mathbf{1}_T/T\| = T^{-1/2}$ , we have that

$$|\hat{b} - b^*| = \mathcal{O}(\epsilon d \log^2(T/p) + T^{-1/2} + T^{-1/2} \|\hat{N} - N^*\|_F).$$

For large enough  $T$ , the term  $T^{-1/2} \|\hat{N} - N^*\|_F$  is a lower-order term, and  $|\hat{b} - b^*| = \mathcal{O}(\epsilon d \log^2(T/p) + T^{-1/2})$ , which is a faster rate than that of  $\|\hat{N} - N^*\|_F$ .  $\square$

## D.1 Persistency of excitation

Below we prove Lemma 1, which claims that  $(f_t)_{t \geq 1}$  satisfies the  $(k, \gamma^2 I, q)$ -BMSB condition. With some additional arguments (Simchowitz et al., 2018; Matni and Tu, 2019), this implies that  $\lambda_{\min}\left(\sum_{t=H}^{T+H-1} f_t f_t^\top\right) = \Omega(\gamma^2 q^2 T)$ , establishing the persistency of excitation.

*Proof of Lemma 1.* Since  $\text{svec}$  is a bijection, every unit vector  $v \in \mathbb{R}^{d_h(d_h+1)/2}$  corresponds to a symmetric matrix  $M \in \mathbb{R}^{d_h \times d_h}$  with unit Frobenius norm. Then,

$$\langle f_{t+i}, v \rangle = \langle \text{svec}(h_{t+i} h_{t+i}^\top), \text{svec}(M) \rangle = h_{t+i}^\top M h_{t+i}.$$

Take  $\Gamma = \gamma^2 I$  for some  $\gamma > 0$ . Then,  $\|v\|_\Gamma = \gamma$ . It suffices to show that for  $i > G$  for some  $G > 0$ ,

$$\mathbb{P}(|h_{t+i}^\top M h_{t+i}| \geq \gamma \mid \mathcal{F}_t) \geq q,$$

since if so, we have

$$\frac{1}{2G} \sum_{i=1}^{2G} \mathbb{P}(|h_{t+i}^\top M h_{t+i}| \geq \gamma \mid \mathcal{F}_t) \geq \frac{1}{2G} \sum_{i=G+1}^{2G} \mathbb{P}(|h_{t+i}^\top M h_{t+i}| \geq \gamma \mid \mathcal{F}_t) \geq q/2,$$

which means  $(f_t)_{t \geq H}$  is  $(2G, \gamma^2 I, q/2)$ -BMSB.

Now let us take a close look at  $h_{t+i} = [y_{(t+i-H+1):(t+i)}; u_{(t+i-H):(t+i-1)}]$ . Since

$$y_{t+i} = C^*(A^*)^i x_t + \sum_{j=1}^i C^*(A^*)^j (B^* u_{t+i-j} + w_{t+i-j}) + v_{t+i},$$

$y_{t+i} \mid \mathcal{F}_t$  is Gaussian with mean  $C^*(A^*)^i x_t$  and covariance determined by  $\sum_{j=1}^i C^*(A^*)^j (B^* u_{t+i-j} + w_{t+i-j}) + v_{t+i}$ , where we note that  $v_{t+i}$  is independent of all other random variables and has full-rank covariance. Hence, for  $i > H$ ,  $h_{t+i} \mid \mathcal{F}_t$  is Gaussian and has full-rank covariance. Then intuitively, since  $\|M\|_F = 1$ ,  $|h_{t+i}^\top M h_{t+i}| \mid \mathcal{F}_t$  is a well-behaved random variable that can exceed some  $\gamma > 0$  with a positive probability  $q$ . Formally, let  $\mu_{t,i} := \mathbb{E}[h_{t+i} \mid \mathcal{F}_t]$ . By Lemma 2, for  $i > H$ , there exists some absolute constant  $a > 0$ , such that

$$\mathbb{E}[|(h_{t+i} - \mu_{t,i})^\top M (h_{t+i} - \mu_{t,i})| \mid \mathcal{F}_t] \geq a \min\{\sigma_u, \sigma_v\} / d_h.$$

By triangle inequality, we have

$$\begin{aligned} |(h_{t+i} - \mu_{t,i})^\top M (h_{t+i} - \mu_{t,i})| &= |h_{t+i}^\top M h_{t+i} + \mu_{t,i}^\top M \mu_{t,i} - 2h_{t+i}^\top M \mu_{t,i}| \\ &\leq |h_{t+i}^\top M h_{t+i}| + |\mu_{t,i}^\top M \mu_{t,i}| + 2|h_{t+i}^\top M \mu_{t,i}|. \end{aligned}$$

Hence,

$$\mathbb{E}[|h_{t+i}^\top M h_{t+i}| \mid \mathcal{F}_t] \geq a \min\{\sigma_u, \sigma_v\} / d_h - \mathbb{E}[|\mu_{t,i}^\top M \mu_{t,i}| + 2|h_{t+i}^\top M \mu_{t,i}| \mid \mathcal{F}_t].$$

Now we argue that for large enough  $i$ ,  $\mathbb{E}[|\mu_{t,i}^\top M \mu_{t,i}| + 2|h_{t+i}^\top M \mu_{t,i}|]$  is negligible. Since matrix  $A^*$  is stable, with probability at least  $1 - p$ ,  $\|x_t\| = \mathcal{O}(d_x^{1/2} \log(T/p))$  for all  $t \geq 0$ . Hence,

$$\|C^*(A^*)^i x_t\| = \mathcal{O}(\alpha(A^*) \rho(A^*)^i d_x^{1/2} \log(T/p)).$$

Then, for  $i > H$ ,

$$\begin{aligned} \mathbb{E}[|\mu_{t,i}^\top M \mu_{t,i}| + 2|h_{t+i}^\top M \mu_{t,i}| \mid \mathcal{F}_t] &= \left\langle \mu_{t,i} \mu_{t,i}^\top, M \right\rangle_F + 2\mathbb{E}[\left\langle \mu_{t,i} h_{t+i}^\top, M \right\rangle_F \mid \mathcal{F}_t] \\ &\leq \|\mu_{t,i} \mu_{t,i}^\top\|_F \cdot \|M\|_F + 2\mathbb{E}[\|\mu_{t,i} h_{t+i}^\top\|_F \cdot \|M\|_F \mid \mathcal{F}_t] \\ &= \|\mu_{t,i}\|^2 + 2\|\mu_{t,i}\| \cdot \mathbb{E}[\|h_{t+i}\| \mid \mathcal{F}_t]. \end{aligned}$$

By definition,  $\mu_{t,i}$  is the concatenation of  $(C^*(A^*)^j x_t)_{i-H+1 \leq j \leq i}$  and zero vectors. Hence,

$$\|\mu_{t,i}\| = \mathcal{O}(d_x^{1/2} d_h^{1/2} \alpha(A^*) \rho(A^*)^i \log(T/p)).$$

Choose  $H \geq \frac{a_1 \log(d_x^{1/2} d_h^2 \alpha(A^*) \log(T/p))}{\log(\rho(A^*)^{-1})}$  for some dimension-free constant  $a_1 > 0$ , such that for  $i > 2H$ , we have

$$\|\mu_{t,i}\|^2 + 2\|\mu_{t,i}\| \cdot \mathbb{E}[\|h_{t+i}\| \mid \mathcal{F}_t] \leq a \min\{\sigma_u, \sigma_v\} / (2d_h).$$

Then, we have the desired lower bound that

$$\mathbb{E}[|h_{t+i}^\top M h_{t+i}| \mid \mathcal{F}_t] \geq a \min\{\sigma_u, \sigma_v\} / (2d_h).$$

On the other hand, since

$$|h_{t+i}^\top M h_{t+i}| = |\langle M, h_{t+i} h_{t+i}^\top \rangle_F| \leq \|M\|_F \|h_{t+i} h_{t+i}^\top\|_F = h_{t+i}^\top h_{t+i},$$

we have  $\mathbb{E}[|h_{t+i}^\top M h_{t+i}|^2 \mid \mathcal{F}_t] \leq \mathbb{E}[\|h_{t+i}\|^4 \mid \mathcal{F}_t]$ . Since  $\|h_{t+i}\| \mid \mathcal{F}_t$  is sub-Gaussian with

$$\|\|h_{t+i}\| \mid \mathcal{F}_t\|_{\psi_2} = \mathcal{O}(\|\mathbb{E}[h_{t+i} h_{t+i}^\top \mid \mathcal{F}_t]\|_2^{1/2}) = \mathcal{O}(1),$$

$\mathbb{E}[|h_{t+i}^\top M h_{t+i}|^2 \mid \mathcal{F}_t] = \mathcal{O}(1)$ . By the Paley-Zygmund inequality, for  $\theta \in [0, 1]$  we have

$$\mathbb{P}(|h_{t+i}^\top M h_{t+i}| \geq \theta a \min\{\sigma_u, \sigma_v\} / (2d_h) \mid \mathcal{F}_t) = \Omega((1 - \theta)^2 a^2 / d_h^2),$$

where the dependence on  $\sigma_u, \sigma_v$  is hidden in  $\Omega(\cdot)$ . By taking  $\theta = 1/2$ , we can see that  $(f_t)_{t \geq H}$  satisfies the  $(k, \gamma^2 I, q)$ -BMSB condition for  $k = 4H$ ,  $\gamma = \Theta(1/d_h)$  and  $q = \Theta(1/d_h^2)$ .  $\square$

## E Linear system identification with noisy measurements

Identifying the time-invariant latent dynamics involves linear regression with time-dependent data and noisy measurements. The following Lemma 8 extends the previous linear system identification result in (Simchowit et al., 2018) to the case with noises in both input and output variables. In Lemma 8,  $\gamma$  and  $q$  are treated as dimension-free constants (in contrast to Lemma 7), which is indeed the case in our application of Lemma 8.

**Lemma 8.** *Let  $(x_t^*)_{t \geq 1}$  be a sequence of  $d_1$ -dimensional Gaussian random vectors adapted to a filtration  $(\mathcal{F}_t)_{t \geq 1}$  with  $\|\mathbb{E}[x_t^* (x_t^*)^\top]\|_2^{1/2} \leq \sigma$ . Define  $y_t^* = A^* x_t^* + e_t$ , where  $A^* \in \mathbb{R}^{d_2 \times d_1}$  and  $e_t \mid \mathcal{F}_t$  is Gaussian with zero mean and  $\|\mathbb{E}[e_t e_t^\top]\|_2^{1/2} \leq \epsilon$ . Define  $y_t = y_t^* + \delta_t^y$  and  $x_t = x_t^* + \delta_t^x$ , where the Gaussian noise vectors  $\delta_t^x$  and  $\delta_t^y$  can be correlated with  $x_t^*$  and  $x_y^*$ , and satisfy  $\|\mathbb{E}[\delta_t^x (\delta_t^x)^\top]\|_2^{1/2} \leq \epsilon_x \leq \sigma$  and  $\|\mathbb{E}[\delta_t^y (\delta_t^y)^\top]\|_2^{1/2} \leq \epsilon_y \leq \sigma$ . Assume that  $(x_t^*)_{t \geq 1}$  satisfies the  $(k, \gamma^2 I, q)$ -BMSB condition and  $\epsilon_x \leq a_0 \gamma^2 q^2 / \sigma$  for some absolute constant  $a_0 > 0$ . Consider*

$$\hat{A} \in \operatorname{argmin}_{A \in \mathbb{R}^{d_2 \times d_1}} \sum_{t=1}^T \|y_t - A x_t\|^2. \quad (\text{E.1})$$

*Then, as long as  $T \geq a_1 k q^{-1} (\log(1/p) + d_1 \log(10/q) + d_1 \log(\sigma \gamma^{-1} d_1 \log(T/p)))$  for some absolute constant  $a_1 > 0$ , we have that with probability at least  $1 - p$ ,*

$$\|\hat{A} - A^*\|_2 = \mathcal{O}((\epsilon_x + \epsilon_y)(d_1 + d_2) \log^2(T/p) + (d_2 + d_1 \log(d_1 \log(T/p)))^{1/2} T^{-1/2}).$$

*Proof.* Let  $X \in \mathbb{R}^{T \times d_1}$  denote the matrix whose  $t$ th row is  $x_t^\top$ . Define  $X^*, Y, E, \Delta_x, \Delta_y$  similarly. To solve the regression problem, we set its gradient to be zero and substitute in  $Y = X^* (A^*)^\top + E + \Delta_y$  to obtain

$$\hat{A} (X^\top X) = A^* (X^*)^\top X + E^\top X + \Delta_y^\top X. \quad (\text{E.2})$$

Substituting in  $X = X^* + \Delta_x$  gives

$$\begin{aligned} (\hat{A} - A^*)((X^*)^\top X^*) &= A^*(X^*)^\top \Delta_x - \hat{A}(\Delta_x^\top \Delta_x + \Delta_x^\top X^* + (X^*)^\top \Delta_x) \\ &\quad + E^\top X^* + E^\top \Delta_x + \Delta_y^\top X^* + \Delta_y^\top \Delta_x. \end{aligned} \quad (\text{E.3})$$

Now we deal with each term on the right-hand side. Since  $(X^*)^\top \Delta_x = \sum_{t=1}^T x_t (\delta_t^x)^\top$ , by triangle inequality,

$$\|(X^*)^\top \Delta_x\|_2 \leq \sum_{t=1}^T \|x_t^* (\delta_t^x)^\top\|_2 \leq \sum_{t=1}^T \|x_t^*\| \cdot \|\delta_t^x\|.$$

Since  $(x_t^*)_{t \geq 1}$  and  $(\delta_t^x)_{t \geq 1}$  are Gaussian, with probability at least  $1 - p$ ,  $\|x_t^*\| = \mathcal{O}(\sigma d_1^{1/2} \log(T/p))$  and  $\|\delta_t^x\| = \mathcal{O}(\epsilon_x d_1^{1/2} \log(T/p))$ . Hence,

$$\|(X^*)^\top \Delta_x\|_2 = \mathcal{O}(\epsilon_x \sigma d_1 T \log^2(T/p)).$$

Similarly, with probability at least  $1 - p$ ,

$$\begin{aligned} \|\Delta_x^\top \Delta_x\|_2 &= \mathcal{O}(\epsilon_x^2 d_1 T \log^2(T/p)), \quad \|E^\top \Delta_x\|_2 = \mathcal{O}(\epsilon \epsilon_x d_1^{1/2} d_2^{1/2} T \log^2(T/p)), \\ \|\Delta_y^\top X^*\|_2 &= \mathcal{O}(\epsilon_y \sigma d_1^{1/2} d_2^{1/2} T \log^2(T/p)), \quad \|\Delta_y^\top \Delta_x\|_2 = \mathcal{O}(\epsilon_x \epsilon_y d_1^{1/2} d_2^{1/2} T \log^2(T/p)). \end{aligned}$$

It remains to bound  $\|\hat{A}\|_2$ . Notice that with probability at least  $1 - p$ ,

$$\|(X^*)^\top X^*\|_2 \leq \sum_{t=1}^T \|x_t\|^2 = \mathcal{O}(\sigma^2 d_1 \log^2(T/p) T).$$

Let  $T_0 := a_1 k q^{-1} (\log(1/p) + d_1 \log(10/q) + d_1 \log(\sigma \gamma^{-1} d_1 \log(T/p)))$  for some absolute constant  $a_1 > 0$ . Then, by (Simchowitz et al., 2018, Appendix D), as long as  $T \geq T_0$ , with probability at least  $1 - p$ ,  $\lambda_{\min}((X^*)^\top X^*) = \gamma^2 q^2 T / 32$ . Since  $X^\top X = (X^*)^\top X^* + \Delta_x^\top X^* + (X^*)^\top \Delta_x + \Delta_x^\top \Delta_x$ ,

$$\lambda_{\min}(X^\top X) \geq \lambda_{\min}((X^*)^\top X^*) - \|\Delta_x^\top X^* + (X^*)^\top \Delta_x + \Delta_x^\top \Delta_x\|_2.$$

Hence, there exists an absolute constant  $a_0 > 0$ , such that as long as  $\epsilon_x \leq a_0 \gamma^2 q^2 / \sigma$ ,  $\lambda_{\min}(X^\top X) = \Omega(\gamma^2 q^2 T)$ , which implies  $\|X^\dagger\|_2 = \mathcal{O}(\gamma^{-1} q^{-1} T^{-1/2})$ . From (E.2), we have

$$\|\hat{A}\|_2 = (\|A^*\|_2 \|X^*\|_2 + \|E\|_2 + \|\Delta_y\|_2) \|X^\dagger\|_2 = \mathcal{O}(\gamma^{-1} q^{-1} (\sigma \|A^*\|_2 + \epsilon + \epsilon_y)) = \mathcal{O}(1),$$

where in the last equality we use  $\epsilon_y \leq \sigma$  and treat  $\gamma, q, \sigma, \epsilon$  as problem-dependent constants.

Finally, by (Simchowitz et al., 2018, Theorem 2.4), as long as  $T \geq T_0$ ,

$$\|E^\top (X^*)^\dagger\|_2 = \mathcal{O}((d_2 + d_1 \log(d_1 \log(T/p)) + \log(1/p))^{1/2} T^{-1/2}).$$

Combining all the above individual bounds with the terms on the right-hand side of (E.3), we have

$$\|\hat{A} - A^*\|_2 = \mathcal{O}((\epsilon_x + \epsilon_y)(d_1 + d_2) \log^2(T/p) + (d_2 + d_1 \log(d_1 \log(T/p)) + \log(1/p))^{1/2} T^{-1/2}),$$

which completes the proof.  $\square$



## F Proof of the main results

In this section, we prove the sample complexity bounds for CoReL and CoReDyL in Theorem 1, which is a simplified version of the following Theorem 2. As we shall see, the proofs for the two algorithms share similar ideas and tools.

**Theorem 2.** *Given an unknown LQG problem satisfying Assumption 1, let  $M^{*l}$  and  $(A^{*l}, B^{*l}, Q^{*l}, R^*)$  be the optimal state representation function and the true system parameters under the normalized parameterization. For a given  $p \in (0, 1)$ , if we run CoReL or CoReDyL for  $T \geq \text{poly}(d_x, d_y, d_u, \log(T/p))$ ,  $H = \Omega\left(\frac{\log(\alpha(\bar{A}^*)(d_y+d_u)^{-1}d_x^{-1/2}T)}{\log(\rho(\bar{A}^*)^{-1})} + \frac{\log(\alpha(A^*)H^2d_x^{1/2}(d_y+d_u)^2\log(T/p))}{\log(\rho(A^*)^{-1})}\right)$ , and  $\sigma_u = \Theta(1)$ , then there exists an orthonormal matrix  $S \in \mathbb{R}^{d_x \times d_x}$ , such that with probability at least  $1 - p$ , the representation function  $\hat{M}$  satisfies*

$$\|\hat{M} - SM^{*l}\|_2 = \mathcal{O}((H(d_y + d_u)d_x)^{15/2}T^{-1/2}\log^6(T/p)),$$

and the suboptimality gap of feedback gain  $\hat{K}$  in system  $(SA^{*l}S^\top, SB^{*l}, SQ^{*l}S^\top, R^*)$  is

$$\mathcal{O}((H(d_y + d_u)d_x)^{15}(d_x \wedge d_u)T^{-1}\log^{12}(T/p)).$$

### F.1 Proof of Theorem 2 for CoReL

*Proof.* Define filtration  $\mathcal{F}_t := \sigma(x_0, v_0, u_0, w_0, v_1, \dots, u_{t-1}, w_{t-1}, v_t)$ . By definition,  $x_t, y_t \in \mathcal{F}_t$ . Hence,  $h_t = [y_{(t-H+1):t}; u_{(t-H):(t-1)}] \in \mathcal{F}_t$  and  $z_t^* \in \mathcal{F}_t$ .

Recall that we define  $f_t := \text{svec}(h_t h_t^\top)$ . By Lemma 1,  $f_t$  is  $(k, \gamma^2 I, q)$ -BMSB for  $k = 4H$ ,  $\gamma = \Theta(1/d_h)$  and  $q = \Theta(1/d_h^2)$ . By Proposition 3 and Lemma 7, there exists some absolute constant  $a_0 > 0$ , such that as long as  $T \geq a_0 H^6 (d_y + d_u)^5 \log(H(d_y + d_u)/p)$ , with probability at least  $1 - p$ ,

$$\|\hat{N} - N^*\|_F = \mathcal{O}((\gamma q)^{-2} d_h E T^{-1/2} \log^4(T/p)) = \mathcal{O}((H(d_y + d_u))^7 d_x^{3/2} T^{-1/2} \log^4(T/p)).$$

By (Tu et al., 2016, Lemma 5.4), there exists an orthonormal matrix  $S$ , such that  $\|\hat{M} - SM^*\|_F$  is on the same order. Since  $\hat{z}_t - Sz_t^* = (\hat{M} - SM^*)h_t - \delta_t$ ,

$$\begin{aligned} \|\text{Cov}(\hat{z}_t - Sz_t^*)\|_2^{1/2} &= \mathcal{O}(\|\text{Cov}((\hat{M} - SM^*)h_t)\|_2^{1/2} + \|\text{Cov}((\bar{A}^*)^H z_{t-H}^*)\|_2^{1/2}) \\ &= \mathcal{O}(\|\mathbb{E}[h_t h_t^\top]\|_2^{1/2} \|\hat{M} - SM^*\|_2 + \|\mathbb{E}[z_{t-H}^* (z_{t-H}^*)^\top]\|_2^{1/2} \|(\bar{A}^*)^H\|_2) \\ &\stackrel{(i)}{=} \mathcal{O}((H(d_y + d_u))^{15/2} d_x^{3/2} T^{-1/2} \log^4(T/p)), \end{aligned}$$

where (i) holds by our choice of  $H$ . Then, by the noisy linear regression bound (Lemma 8), for  $T$  greater than a constant polynomial in the problem parameters, we have

$$\begin{aligned} \|[\hat{A}, \hat{B}] - S[A^*, B^*]\|_2 &= \mathcal{O}((H(d_y + d_u))^{15/2} d_x^{3/2} T^{-1/2} \log^4(T/p) \cdot d_x \log^2(T/p)) \\ &= \mathcal{O}((H(d_y + d_u))^{15/2} d_x^{5/2} T^{-1/2} \log^6(T/p)). \end{aligned}$$

By the noisy quadratic regression bound (Lemma 7) and  $\text{svec}(z_t^*(z_t^*)^\top)$  being  $(\Theta(1), \Theta(1/d_x), \Theta(1/d_x^2))$ -BMSM,

$$\begin{aligned}\|\hat{Q} - SQ^*S^\top\|_F &= \mathcal{O}((H(d_y + d_u))^{15/2}d_x^{3/2}T^{-1/2}\log^4(T/p) \cdot d_x^3d_x\log^2(T/p) + d_x^6d_xd_x^{1/2}T^{-1/2}\log^4(T/p)) \\ &= \mathcal{O}((H(d_y + d_u)d_x)^{15/2}T^{-1/2}\log^6(T/p)).\end{aligned}$$

Hence,  $\|\hat{A} - SA^*S^\top\|_2$ ,  $\|\hat{B} - SB^*S^\top\|_2$  and  $\|\hat{Q} - SQ^*S^\top\|_2$  are all bounded by  $\mathcal{O}((H(d_y + d_u)d_x)^{15/2}T^{-1/2}\log^6(T/p))$ . By (Mania et al., 2019), for  $T$  greater than a constant polynomial in the problem parameters,

$$\|\hat{K} - K^*S^\top\|_2 = \mathcal{O}((H(d_y + d_u)d_x)^{15/2}T^{-1/2}\log^6(T/p))$$

is on the same order, and  $\hat{K}$  has a suboptimality gap of

$$\mathcal{O}((H(d_y + d_u)d_x)^{15}(d_x \wedge d_u)T^{-1}\log^{12}(d_xT/p)).$$

□

## F.2 Proof of Theorem 2 for CoReDyL

*Proof.* Define filtration  $\mathcal{F}_t := \sigma(x_0, v_0, u_0, w_0, v_1, \dots, u_{t-1}, w_{t-1}, v_t)$ . By definition,  $x_t, y_t \in \mathcal{F}_t$ . Hence,  $h_t = [y_{(t-H+1):t}; u_{(t-H):(t-1)}] \in \mathcal{F}_t$  and  $z_t^* \in \mathcal{F}_t$ .

Recall that we define  $f_t := \text{svec}(h_t h_t^\top)$ . By Lemma 1,  $f_t$  is  $(k, \gamma^2 I, q)$ -BMSB for  $k = 4H$ ,  $\gamma = \Theta(1/d_h)$  and  $q = \Theta(1/d_h^2)$ . By Proposition 3 and Lemma 7, there exists some absolute constant  $a_0 > 0$ , such that as long as  $T \geq a_0 H^6 (d_y + d_u)^5 \log(H(d_y + d_u)/p)$ , with probability at least  $1 - p$ ,

$$\begin{aligned}\|\hat{N} - N^*\|_F &= \mathcal{O}((H(d_y + d_u))^7 d_x^{3/2} T^{-1/2} \log^4(T/p)), \\ \|\hat{N}_1 - N_1^*\|_F &= \mathcal{O}((H(d_y + d_u) + d_x)^7 d_x^{3/2} T^{-1/2} \log^4(T/p)).\end{aligned}$$

By (Tu et al., 2016, Lemma 5.4), there exists orthonormal matrices  $S, S_1$ , such that  $\|\hat{M} - SM^*\|_F$  and  $\|\hat{M}_1 - S_1 M_1^*\|_F$  are on the same order of  $\|\hat{N} - N^*\|_F$  and  $\|\hat{N}_1 - N_1^*\|_F$ , respectively. Then,

$$\|\tilde{B} - S_1 B^*\|_2 = \mathcal{O}((H(d_y + d_u) + d_x)^7 d_x^{3/2} T^{-1/2} \log^4(T/p)).$$

Since  $[S_1 A^* S^\top S M^*, S_1 B^*] = S_1 M_1^*$ , by the perturbation bounds of the Moore-Penrose inverse (Wedin, 1973), we have

$$\|\tilde{A} - S_1 A^* S^\top\|_2 = \mathcal{O}((H(d_y + d_u) + d_x)^7 d_x^{3/2} T^{-1/2} \log^4(T/p)).$$

To align  $\tilde{A}$  with  $SA^*S^\top$ , we compute another matrix  $\hat{S}_0$  by solving

$$[\hat{z}_{H+1}, \dots, \hat{z}_{T+H}] = \hat{S}_0 \hat{M}_1 [[h_H; u_H], \dots, [h_{T+H-1}; u_{T+H-1}]].$$

Since  $\hat{z}_t - Sz_t^* = (\hat{M} - SM^*)h_t - \delta_t$ ,

$$\begin{aligned}\|\text{Cov}(\hat{z}_t - Sz_t^*)\|_2^{1/2} &= \mathcal{O}(\|\text{Cov}((\hat{M} - SM^*)h_t)\|_2^{1/2} + \|\text{Cov}((\bar{A}^*)^H z_{t-H}^*)\|_2^{1/2}) \\ &= \mathcal{O}(\|\mathbb{E}[h_t h_t^\top]\|_2^{1/2} \|\hat{M} - SM^*\|_2 + \|\mathbb{E}[z_{t-H}^* (z_{t-H}^*)^\top]\|_2^{1/2} \|(\bar{A}^*)^H\|_2) \\ &\stackrel{(i)}{=} \mathcal{O}((H(d_y + d_u))^{15/2} d_x^{3/2} T^{-1/2} \log^4(T/p)),\end{aligned}$$

where (i) holds by our choice of  $H$ . Similarly,

$$\|\text{Cov}(\hat{M}_1[h_t; u_t] - S_1 z_{t+1}^*)\|_2^{1/2} = \mathcal{O}((H(d_y + d_u) + d_x)^{15/2} d_x^{3/2} T^{-1/2} \log^4(T/p)).$$

Since  $S[z_{H+1}^*, \dots, z_{H+T}^*] = SS_1^\top S_1[z_{H+1}^*, \dots, z_{H+T-1}^*]$ , by the noisy linear regression bound (Lemma 8), for  $T$  greater than a constant polynomial in the problem parameters, we have

$$\begin{aligned} \|\hat{S}_0 - SS_1^\top\|_2 &= \mathcal{O}((H(d_y + d_u) + d_x)^{15/2} d_x^{3/2} T^{-1/2} \log^4(T/p) \cdot d_x \log^2(T/p)) \\ &= \mathcal{O}((H(d_y + d_u) + d_x)^{15/2} d_x^{5/2} T^{-1/2} \log^6(T/p)). \end{aligned}$$

Hence,

$$\begin{aligned} \|\hat{A} - SA^*S^\top\|_2 &= \|\hat{S}_0 \tilde{A} - SS_1^\top S_1 A^*S^\top\|_2 = \|(\hat{S}_0 - SS_1^\top) \tilde{A}\|_2 + \|SS_1^\top (\tilde{A} - S_1 A^*S^\top)\|_2 \\ &= \mathcal{O}((H(d_y + d_u) + d_x)^{15/2} d_x^{5/2} T^{-1/2} \log^6(T/p)), \end{aligned}$$

and  $\|\hat{B} - SB^*\|_2$  has the same order. By the noisy quadratic regression bound (Lemma 7) and  $\text{svec}(z_t^* (z_t^*)^\top)$  being  $(\Theta(1), \Theta(1/d_x), \Theta(1/d_x^2))$ -BMSM,

$$\begin{aligned} \|\hat{Q} - SQ^*S^\top\|_F &= \mathcal{O}((H(d_y + d_u))^{15/2} d_x^{3/2} T^{-1/2} \log^4(T/p) \cdot d_x^3 d_x \log^2(T/p) + d_x^6 d_x d_x^{1/2} T^{-1/2} \log^4(T/p)) \\ &= \mathcal{O}((H(d_y + d_u) d_x)^{15/2} T^{-1/2} \log^6(T/p)). \end{aligned}$$

Hence,  $\|\hat{A} - SA^*S^\top\|_2$ ,  $\|\hat{B} - SB^*S^\top\|_2$  and  $\|\hat{Q} - SQ^*S^\top\|_2$  are all bounded by  $\mathcal{O}((H(d_y + d_u) d_x)^{15/2} T^{-1/2} \log^6(T/p))$ . By (Mania et al., 2019), for  $T$  greater than a constant polynomial in the problem parameters,

$$\|\hat{K} - K^*S^\top\|_2 = \mathcal{O}((H(d_y + d_u) d_x)^{15/2} T^{-1/2} \log^6(T/p))$$

is on the same order, and  $\hat{K}$  has a suboptimality gap of

$$\mathcal{O}((H(d_y + d_u) d_x)^{15} (d_x \wedge d_u) T^{-1} \log^{12}(d_x T/p)),$$

which completes the proof.  $\square$