# Revisiting Ho-Kalman based system identification: robustness and finite-sample analysis

Samet Oymak, *Member, IEEE*    and    Necmiye Ozay, *Senior Member, IEEE*

*Abstract*— **We consider the problem of learning a realization for a linear time-invariant (LTI) dynamical system from input/output data. Given a single input/output trajectory, we provide finite time analysis for learning the system's Markov parameters, from which a balanced realization is estimated using the classical Ho-Kalman algorithm. By proving a robustness result for the Ho-Kalman algorithm and combining it with the sample complexity results for Markov parameters, we show how much data is needed to approximate the balanced realization of the system up to a desired accuracy with high probability.**

*Index Terms*— **system identification, sample complexity, balanced realization, Markov parameters**

## I. INTRODUCTION

Many modern control design techniques rely on the existence of a fairly accurate state-space model of the plant to be controlled. Although in some cases a model can be obtained from first principles, there are many situations in which a model should be learned from input/output data. Classical results in system identification provide asymptotic convergence guarantees for learning models from data [1], [2], [3]. Finite sample complexity properties have also been discussed in system identification literature [4], [5], [6], [7], with various different types of assumptions (cf. [8] for a recent survey); however earlier results rely on assumptions that might not hold for input/output models or are conservative at times [9].

There is recent interest from the machine learning community in data-driven control and non-asymptotic analysis. Putting aside the reinforcement learning literature and restricting our attention to linear state-space models, the work in this area can be divided into two categories: (i) directly learning the control inputs to optimize a control objective or analyzing the predictive power of the learned representation [10], [11], [12], (ii) learning the parameters of the system model from limited data [13], [14], [15], [16], [9], [17]. For the former problem, the focus has been on exploration/exploitation type formulations and regret analysis. Since the goal is to learn how to control the system to achieve a specific task, the system is not necessarily fully learned. On the other hand, the latter problem aims to learn a general purpose model that can be used in different control tasks, for instance, by combining it

with robust control techniques [15], [18], [16]. The focus for the latter work has been to analyze data–accuracy trade-offs.

In this paper we focus on learning a realization for an LTI system from a single *input/output* trajectory. This setting is significantly more challenging than earlier studies that assume that (multiple independent) *state* trajectories are available [15], [9], [19], [20], [21]. One of our main contributions is to derive sample complexity results in learning the Markov parameters, to be precisely defined later, of the system using a least squares algorithm [22]. Markov parameters play a central role in system identification [1] and they can also be directly used in control design when the system model itself is not available [23], [24], [25]. When only input/output data is available, it is well known that the system matrices can be identified only up to a similarity transformation even in the noise-free case but Markov parameters are identifiable. Therefore, we focus on obtaining a realization. One classical technique to derive a realization from the Markov parameters is the Ho-Kalman (a.k.a., eigensystem realization algorithm – ERA) algorithm [26]. The Ho-Kalman algorithm constructs a balanced realization for the system from the singular value decomposition of the Hankel matrix of the Markov parameters (see, e.g., [27], [28].). By proving a robustness result for the Ho-Kalman algorithm, i.e., small changes in the inputs of the algorithm lead to small changes in its outputs, and combining it with the sample complexity results, we show how much data is needed to learn a balanced realization of the system up to a desired accuracy with high probability.

Most common analysis for system identification algorithms is that of consistency [1]. Other asymptotic analysis include that of asymptotic normality, which shows that the errors in the estimates of system matrices follow a normal distribution [2], [29] as the number of samples goes to infinity. On the other hand, non-asymptotic, finite sample results in system identification (cf. [8]) focus on numerical evaluation of probabilistic bounds on the accuracy of the estimates as a function of the number of samples. Our analysis is different in that it reveals the dependency of the accuracy of the estimates on the number of samples and on the system theoretic properties of the underlying system such as norms of Gramians, spectral radius, etc. In this sense, our results are useful in understanding what systems are easier to identify with the Ho-Kalman algorithm. Moreover, our bounds can give theoretical insights into practical heuristics that are used for establishing high-probability confidence bounds, such as bootstrapping [15].

S. Oymak is with the Department of Electrical and Computer Engineering, University of California, Riverside. N. Ozay is with the Department of Electrical Engineering and Computer Science, University of Michigan, Ann Arbor.

A short version of this work appeared in [30] where we provide preliminary guarantees for the system identification problem. This work strengthens the results of [30] and also provides the necessary technical framework and the associated proofs. Two important changes from [30] are as follows. First, [30] uses a bound on the spectral norm of the state matrix as the stability assumption on the system. In contrast, this manuscript uses spectral radius which is much less restrictive and more conventional in the control literature. Secondly, in Section IV, we provide new bounds on Hankel matrix estimation by leveraging the stability of the system. We remark that after this manuscript appeared on arXiv, there have been other interesting works on the finite sample system identification problem, some following up on the Ho-Kalman analysis framework we introduced. Sarkar et al. [31] provide non-asymptotic guarantees for systems with unknown orders, where they directly estimate the Hankel matrix of the system which leads to a quadratic growth in sample complexity as a function of model order. This is in contrast to our two stage approach which first estimates the Markov parameters. Simchowitz et al. [32] consider a semi-parametric noise model and studies prefiltered least-squares. Mania et al. [33] use finite-sample estimates of system dynamics to provide performance bound for Kalman filter. Tsiamis and Pappas [34] study stochastic systems which are purely driven by noise.

## II. PROBLEM SETUP

We first introduce the basic notation. Spectral norm $\|\cdot\|$ returns the largest singular value of a matrix. Multivariate normal distribution with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$ is denoted by $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. $\boldsymbol{X}^*$ denotes the transpose of a matrix $\boldsymbol{X}$. $\boldsymbol{X}^\dagger$ returns the Moore–Penrose inverse of the matrix $\boldsymbol{X}$. Covariance matrix of a random vector $\boldsymbol{v}$ is denoted by $\boldsymbol{\Sigma}(\boldsymbol{v})$. $\mathrm{tr}(\cdot)$ returns the trace of a matrix. $c, C, c_0, c_1, \dots$ stand for absolute constants. $\gtrsim, \lesssim$ denote inequalities that hold up to an absolute constant.

Suppose we have an observable and controllable linear system characterized by the system matrices $\boldsymbol{A} \in \mathbb{R}^{n \times n}, \boldsymbol{B} \in \mathbb{R}^{n \times p}, \boldsymbol{C} \in \mathbb{R}^{m \times n}, \boldsymbol{D} \in \mathbb{R}^{m \times p}$ and this system evolves according to

$$\boldsymbol{x}_{t+1} = \boldsymbol{A}\boldsymbol{x}_t + \boldsymbol{B}\boldsymbol{u}_t + \boldsymbol{w}_t, \quad (\mathrm{II.1})$$
$$\boldsymbol{y}_t = \boldsymbol{C}\boldsymbol{x}_t + \boldsymbol{D}\boldsymbol{u}_t + \boldsymbol{z}_t. \quad (\mathrm{II.2})$$

Our goal is to learn the characteristics of this system and to provide finite sample bounds on the estimation accuracy. Given a horizon $T$, we will learn the first $T$ Markov parameters of the system. The first Markov parameter is the matrix $\boldsymbol{D}$, and the remaining parameters are the set of matrices $\{\boldsymbol{C}\boldsymbol{A}^i\boldsymbol{B}\}_{i=0}^{T-2}$. As it will be discussed later on, by learning these parameters,

- we can provide bounds on how well $\boldsymbol{y}_t$ can be estimated for a future time $t$,
- we can identify the state-space matrices $\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}, \boldsymbol{D}$ (up to a similarity transformation).

**Problem setup:** We assume that $\{\boldsymbol{u}_t, \boldsymbol{w}_t, \boldsymbol{z}_t\}_{t=1}^\infty$ are vectors that are independent of each other with distributions $\boldsymbol{u}_t \sim$

$\mathcal{N}(0, \sigma_u^2 \boldsymbol{I}_p)$, $\boldsymbol{w}_t \sim \mathcal{N}(0, \sigma_w^2 \boldsymbol{I}_n)$, and $\boldsymbol{z}_t \sim \mathcal{N}(0, \sigma_z^2 \boldsymbol{I}_m)$[1]. $\boldsymbol{u}_t$ is the input vector which is known to us. $\boldsymbol{w}_t$ and $\boldsymbol{z}_t$ are the process and measurement noise vectors respectively. We also assume that the initial condition of the hidden state is $\boldsymbol{x}_1 = 0$. Observe that Markov parameters can be found if we have access to cross correlations $\mathbb{E}[\boldsymbol{y}_t \boldsymbol{u}_{t-k}^*]$. In particular, we have the identities

$$\mathbb{E}\left[\frac{\boldsymbol{y}_t \boldsymbol{u}_{t-k}^*}{\sigma_u^2}\right] = \begin{cases} \boldsymbol{D} & \text{if} \quad k = 0, \\ \boldsymbol{C}\boldsymbol{A}^{k-1}\boldsymbol{B} & \text{if} \quad k \geq 1 \end{cases}.$$

Hence, if we had access to infinitely many independent $(\boldsymbol{y}_t, \boldsymbol{u}_{t-k})$ pairs, our task could be accomplished by a simple averaging. In this work, we will show that, one can robustly learn these matrices from a small amount of data generated from a single realization of the system trajectory. The challenge is efficiently using finite and dependent data points to perform reliable estimation. Observe that, our problem is identical to learning the concatenated matrix $\boldsymbol{G}$ defined as

$$\boldsymbol{G} = [\boldsymbol{D}, \ \boldsymbol{C}\boldsymbol{B}, \ \boldsymbol{C}\boldsymbol{A}\boldsymbol{B}, \ \dots, \ \boldsymbol{C}\boldsymbol{A}^{T-2}\boldsymbol{B}] \in \mathbb{R}^{m \times Tp}.$$

Next section describes our input and output data. Based on this, we formulate a least-squares procedure that estimates $\boldsymbol{G}$. The estimate $\hat{\boldsymbol{G}}$ will play a critical role in the identification of the system matrices.

### A. Least-Squares Procedure

To describe the estimation procedure, we start by explaining the data collection process. Given a single input/output trajectory $\{\boldsymbol{y}_t, \boldsymbol{u}_t\}_{t=1}^{\bar{N}}$, we generate $N$ subsequences of length $T$, where $\bar{N} = T + N - 1$ and $N \geq 1$. To ease representation, we organize the data $\boldsymbol{u}_t$ and the noise $\boldsymbol{w}_t$ into length $T$ chunks denoted by the following vectors,

$$\bar{\boldsymbol{u}}_t = [\boldsymbol{u}_t^* \ \boldsymbol{u}_{t-1}^* \ \dots \ \boldsymbol{u}_{t-T+1}^*]^* \in \mathbb{R}^{Tp}, \quad (\mathrm{II.3})$$
$$\bar{\boldsymbol{w}}_t = [\boldsymbol{w}_t^* \ \boldsymbol{w}_{t-1}^* \ \dots \ \boldsymbol{w}_{t-T+1}^*]^* \in \mathbb{R}^{Tn}. \quad (\mathrm{II.4})$$

In a similar fashion to $\boldsymbol{G}$ define the matrix,

$$\boldsymbol{F} = [\boldsymbol{0} \ \boldsymbol{C} \ \boldsymbol{C}\boldsymbol{A} \ \dots \ \boldsymbol{C}\boldsymbol{A}^{T-2}] \in \mathbb{R}^{m \times Tn}.$$

To establish an explicit connection to Markov parameters, $\boldsymbol{y}_t$ can be expanded recursively until $t - T + 1$ to relate the output to the input $\bar{\boldsymbol{u}}_t$ and Markov parameter matrix $\boldsymbol{G}$ as follows,

$$\begin{aligned} \boldsymbol{y}_t &= \boldsymbol{C}\boldsymbol{x}_t + \boldsymbol{D}\boldsymbol{u}_t + \boldsymbol{z}_t, \\ &= \boldsymbol{C}(\boldsymbol{A}\boldsymbol{x}_{t-1} + \boldsymbol{B}\boldsymbol{u}_{t-1} + \boldsymbol{w}_{t-1}) + \boldsymbol{D}\boldsymbol{u}_t + \boldsymbol{z}_t, \\ &= \boldsymbol{G}\bar{\boldsymbol{u}}_t + \boldsymbol{F}\bar{\boldsymbol{w}}_t + \boldsymbol{z}_t + \boldsymbol{e}_t, \quad (\mathrm{II.5}) \end{aligned}$$

where, $\boldsymbol{e}_t = \boldsymbol{C}\boldsymbol{A}^{T-1}\boldsymbol{x}_{t-T+1}$ corresponds to the error due to the effect of the state at time $t - T + 1$. With this relation, we will use $(\bar{\boldsymbol{u}}_t, \boldsymbol{y}_t)_{t=T}^{\bar{N}}$ as inputs and outputs of our regression problem. We treat $\bar{\boldsymbol{w}}_t$, $\boldsymbol{z}_t$, and $\boldsymbol{e}_t$ as additive noise and attempt to estimate $\boldsymbol{G}$ from covariates $\bar{\boldsymbol{u}}_t$. Note that, the noise terms are zero-mean including $\boldsymbol{e}_t$ since we assumed $\boldsymbol{x}_1 = 0$. With these in mind, we form the following least-squares problem,

$$\hat{\boldsymbol{G}} = \underset{\boldsymbol{X} \in \mathbb{R}^{m \times Tp}}{\arg\min} \sum_{t=T}^{\bar{N}} \|\boldsymbol{y}_t - \boldsymbol{X}\bar{\boldsymbol{u}}_t\|_{\ell_2}^2.$$

[1]While we assume diagonal covariance throughout the paper, we believe our proof strategy can be adapted to arbitrary covariance matrices.

Defining our label matrix $\boldsymbol{Y}$ and input data matrix $\boldsymbol{U}$ as,

$$\boldsymbol{Y} = [\boldsymbol{y}_T, \ \boldsymbol{y}_{T+1}, \ \ldots, \ \boldsymbol{y}_{\bar{N}}]^* \in \mathbb{R}^{N \times m} \quad \text{and} \quad \text{(II.6)}$$
$$\boldsymbol{U} = [\bar{\boldsymbol{u}}_T, \ \bar{\boldsymbol{u}}_{T+1}, \ \ldots, \ \bar{\boldsymbol{u}}_{\bar{N}}]^* \in \mathbb{R}^{N \times Tp},$$

we obtain the minimization $\min_{\boldsymbol{X}} \|\boldsymbol{Y} - \boldsymbol{U}\boldsymbol{X}^*\|_F^2$. Hence, the least-squares solution $\hat{\boldsymbol{G}}$ is given by

$$\hat{\boldsymbol{G}} = (\boldsymbol{U}^\dagger \boldsymbol{Y})^*, \quad \text{(II.7)}$$

where $\boldsymbol{U}^\dagger = (\boldsymbol{U}^*\boldsymbol{U})^{-1}\boldsymbol{U}^*$ is the left pseudo-inverse of $\boldsymbol{U}$. Ideally, we would like the estimation error $\|\boldsymbol{G} - \hat{\boldsymbol{G}}\|_F^2$ to be small. Our main result bounds the norm of the error as a function of the sample size $N$ and noise levels $\sigma_w$ and $\sigma_z$.

## III. RESULTS ON LEARNING MARKOV PARAMETERS

Let $\rho(\cdot)$ denote the spectral radius of a matrix which is the largest absolute value of its eigenvalues. Our results in this section apply to stable systems where $\rho(\boldsymbol{A}) < 1$. Additionally we need a related quantity involving $\boldsymbol{A}$ which is the ratio between the exponents of the spectral norm and the square-root of the spectral radius defined as

$$\Phi(\boldsymbol{A}) = \sup_{\tau \geq 0} \frac{\|\boldsymbol{A}^\tau\|}{\rho(\boldsymbol{A})^{\tau/2}}.$$

$\Phi(\boldsymbol{A})$ is guaranteed to be finite thanks to Gelfand's formula which states that $\sup_{\tau \geq 0} \|\boldsymbol{A}^\tau\|\rho^{-\tau}$ is finite if $\rho > \rho(\boldsymbol{A})$. In our case, we chose $1 > \rho = \sqrt{\rho(\boldsymbol{A})} > \rho(\boldsymbol{A})$. Another important parameter is the steady state covariance matrix of $\boldsymbol{x}_t$ which is given by

$$\boldsymbol{\Gamma}_\infty = \sum_{i=0}^{\infty} \sigma_w^2 \boldsymbol{A}^i (\boldsymbol{A}^*)^i + \sigma_u^2 \boldsymbol{A}^i \boldsymbol{B} \boldsymbol{B}^* (\boldsymbol{A}^*)^i.$$

This is essentially the sum of scaled controllability Gramians with respect to the process noise and to the control input. It is rather trivial to show that for all $t \geq 1$, $\boldsymbol{\Sigma}(\boldsymbol{x}_t) \preceq \boldsymbol{\Gamma}_\infty$. We will use $\boldsymbol{\Gamma}_\infty$ to bound the error $\boldsymbol{e}_t$ due to the unknown state at time $t - T + 1$. Following the definition of $\boldsymbol{e}_t$, we have that $\|\boldsymbol{\Sigma}(\boldsymbol{e}_t)\| \leq \|\boldsymbol{C}\boldsymbol{A}^{T-1}\|^2 \|\boldsymbol{\Gamma}_\infty\|$. We characterize the impact of $\boldsymbol{e}_t$ by its "effective standard deviation" $\sigma_e$ that is obtained by scaling the bound on $\sqrt{\|\boldsymbol{\Sigma}(\boldsymbol{e}_t)\|}$ by an additional factor $\Phi(\boldsymbol{A})\sqrt{T/(1-\rho(\boldsymbol{A})^T)}$ which yields,

$$\sigma_e = \Phi(\boldsymbol{A})\|\boldsymbol{C}\boldsymbol{A}^{T-1}\|\sqrt{\frac{T\|\boldsymbol{\Gamma}_\infty\|}{1-\rho(\boldsymbol{A})^T}}. \quad \text{(III.1)}$$

Throughout, we assume $Np, \ Tp \geq 2$ which helps simplify the notation. Before stating our main result in Theorem 3.2, we present a simplified version that captures the problem dependencies in terms of the *total standard deviations* $\sigma_z + \sigma_e + \sigma_w\|\boldsymbol{F}\|$ and the *total dimension* $q = m + p + n$.

*Theorem 3.1:* Set $q = p+n+m$. Suppose $\rho(\boldsymbol{A})^T \leq 0.99$ and for a proper constant $c >$, the sample size obeys

$$\frac{N}{\log^2(Nq)} \geq N_0 := cTq\log^2(Tq). \quad \text{(III.2)}$$

Given observations of a single trajectory until time $\bar{N} = N + T - 1$, with high probability[2], the least-square estimator of the Markov parameter matrix obeys

$$\|\hat{\boldsymbol{G}} - \boldsymbol{G}\| \leq \frac{\sigma_z + \sigma_e + \sigma_w\|\boldsymbol{F}\|\log(Nq)}{\sigma_u}\sqrt{\frac{N_0}{N}}.$$

We first describe the basic proof idea. Following equation (II.5), to further simplify the notation, define the matrices

$$\boldsymbol{W} = [\bar{\boldsymbol{w}}_T, \ \bar{\boldsymbol{w}}_{T+1}, \ \ldots, \ \bar{\boldsymbol{w}}_{\bar{N}}]^* \in \mathbb{R}^{N \times Tn}, \quad \text{(III.3)}$$
$$\boldsymbol{E} = [\boldsymbol{e}_T, \ \boldsymbol{e}_{T+1}, \ \ldots, \ \boldsymbol{e}_{\bar{N}}]^* \in \mathbb{R}^{N \times n},$$
$$\boldsymbol{Z} = [\boldsymbol{z}_T, \ \boldsymbol{z}_{T+1}, \ \ldots, \ \boldsymbol{z}_{\bar{N}}]^* \in \mathbb{R}^{N \times m}.$$

With these variables, we have the system of equations

$$\boldsymbol{Y} = \boldsymbol{U}\boldsymbol{G}^* + \boldsymbol{E} + \boldsymbol{Z} + \boldsymbol{W}\boldsymbol{F}^*.$$

Following (II.7), estimation error is given by

$$(\hat{\boldsymbol{G}} - \boldsymbol{G})^* = (\boldsymbol{U}^*\boldsymbol{U})^{-1}\boldsymbol{U}^*(\boldsymbol{W}\boldsymbol{F}^* + \boldsymbol{Z} + \boldsymbol{E}). \quad \text{(III.4)}$$

Hence, the spectral norm of the error can be bounded as

$$\|(\hat{\boldsymbol{G}} - \boldsymbol{G})^*\| \leq \|(\boldsymbol{U}^*\boldsymbol{U})^{-1}\|(\|\boldsymbol{U}^*\boldsymbol{W}\|\|\boldsymbol{F}^*\| + \|\boldsymbol{U}^*\boldsymbol{Z}\| + \|\boldsymbol{U}^*\boldsymbol{E}\|). \quad \text{(III.5)}$$

The proofs of both this theorem and the next follow by individually bounding each of the terms appearing in the above bound. The bounds on $\|(\boldsymbol{U}^*\boldsymbol{U})^{-1}\|$ and $\|\boldsymbol{U}^*\boldsymbol{W}\|$ are obtained by using the properties of random circulant matrices in Appendix III. $\|\boldsymbol{U}^*\boldsymbol{Z}\|$ is arguably the simplest term due to $\boldsymbol{Z}$ being an i.i.d. Gaussian matrix. It is bounded via Lemma 1.1. Finally, $\|\boldsymbol{U}^*\boldsymbol{E}\|$ term is addressed by employing a martingale based argument in Appendix IV.

**Remark:** Our result is stated in terms of the spectral norm error $\|\hat{\boldsymbol{G}} - \boldsymbol{G}\|$. One can deduce the following Frobenius norm bound by naively bounding $\sigma_e, \sigma_z$ terms and swapping $\|\boldsymbol{F}\|$ term by $\|\boldsymbol{F}\|_F$ following Eqs. (III.4) and (III.5). This yields, $\|\hat{\boldsymbol{G}} - \boldsymbol{G}\|_F \leq \frac{(\sigma_z + \sigma_e)\sqrt{m} + \sigma_w\|\boldsymbol{F}\|_F\log(Nq)}{\sigma_u}\sqrt{\frac{N_0}{N}}$.

Our bound individually accounts for the the process noise sequence $\{\boldsymbol{w}_\tau\}_{\tau=t-T+1}^t$, measurement noise $\boldsymbol{z}_t$, and the contribution of the unknown state $\boldsymbol{x}_{t-T+1}$. At a high-level the result indicates that the accuracy of Markov parameter estimates increases as the variance of the input increases and the effective noise or the spectral radius of the $A$ matrix decreases. There is a $\|\boldsymbol{C}\boldsymbol{A}^{T-1}\|$ multiplier inside the unknown state component $\sigma_e$ hence larger $T$ implies smaller $\sigma_e$. On the other hand, larger $T$ increases the size of the $\boldsymbol{G}$ matrix as its dimensions are $m \times Tp$. As a result sample size should grow proportional to $T$ which is reflected within the sample complexity term $N_0$ which grows proportional to $T$ (ignoring $\log$ terms). This highlights that, when $N$ is fixed, there is a sweet spot for the choice of $T$ as it should be large enough to ensure small $\sigma_e$ but as small enough to satisfy the requirement (III.2). In the subsequent discussion, Theorem 4.2 provides such a $T$ choice which ensures the $\sigma_e$ term becomes sufficiently small.

Observe than $Tp$ is the minimum observation period for estimating $\boldsymbol{G}$ since there are $mTp$ unknowns and we get to observe $m$ measurements at each time step. Hence, even if

---

[2]Precise statement on the probability of success is provided in the proof.

one was solving our regression problem with an ideal input data $U$ (e.g. with independent standardized entries) - the result would still require at least $N \gtrsim Tp$ samples and the estimation error rate would decay as $\sqrt{1/N}$ [35]. This is because the input matrix $U$ becomes tall as soon as $N \geq Tp$ at which point the system of equations becomes invertible. Hence ideally (II.7) should work as soon as $N \gtrsim Tp$. Instead Theorem 3.1 is operational in the regime $N \gtrsim T(p+n+m)$ up to log factors. For systems where the state dimension $n$ is much larger than the number of sensors $m$ and input dimension $p$, Theorem 3.1 can be suboptimal. Our main theorem, stated below, is a refined version of Theorem 3.1 and **achieves the optimal sample complexity** of $N \gtrsim Tp$ (up to log factors) which is independent of $n$. This theorem also carefully quantifies the contribution of each noise type to the overall estimation error.

*Theorem 3.2:* Suppose system is stable (i.e. $\rho(A) < 1$) and $\frac{N}{\log^2(Np)} \geq cTp\log^2(Tp)$. We observe a trajectory until time $\bar{N} = N + T - 1$. Set $q = p + n$ and $N_w = cTq\log^2(Tq)\log^2(Nq)$. Then, with high probability, the least-square estimator of the Markov parameter matrix obeys

$$\|\hat{G} - G\| \leq \frac{R_w + R_e + R_z}{\sigma_u \sqrt{N}}, \qquad \text{(III.6)}$$

where $R_w, R_e, R_z$ are given by

$$R_z = 8\sigma_z \sqrt{Tp + m},$$
$$R_w = \sigma_w \|F\| (\sqrt{N_w} + N_w/\sqrt{N}),$$
$$R_e = C\sigma_e \sqrt{\left(1 + \frac{mT}{N(1 - \rho(A)^{T/2})}\right)(Tp + m)}.$$

One can obtain Theorem 3.1 from Theorem 3.2 as follows. When $N \geq N_0 \log^2(Nq) \geq N_w$, $R_w$ satisfies $R_w \leq \sigma_w \|F\| \sqrt{N_w} \leq \sigma_w \|F\| \sqrt{N_0} \log(Nq)$. Similarly, when $\rho(A)^T$ is bounded away from 1 by a constant and $N \geq N_0 \geq c \cdot Tm$, $R_e$ satisfies $R_e \leq 2C\sigma_e \sqrt{Tp + m} \leq \sigma_e \sqrt{N_0}$.

A key advantage of Theorem 3.2 is that it applies in the regime $Tp \lesssim N \lesssim T(p+n+m)$. Additionally, Theorem 3.2 provides tighter individual error bounds for the $\sigma_z, \sigma_w, \sigma_e$ terms and explicitly characterizes the dependence on $\rho(A)$ inside the $R_e$ term.

Theorem 3.2 can be improved in a few directions. Some of the log factors that appear in our sample size might be spurious. These terms are arising from a theorem borrowed from Krahmer et al. [36]; which actually has a stronger implication than what we need in this work. We also believe (III.1) is overestimating the correct dependence by a factor of $\sqrt{T}$.

### A. Predicting the System Output via Markov Parameters

The following lemma illustrates how learning Markov parameters helps us bound the prediction error.

*Lemma 3.3 (Predicting $y_T$):* Suppose $x_1 = 0$ and $z_t \sim \mathcal{N}(0, \sigma_z^2 I)$, $u_t \sim \mathcal{N}(0, \sigma_u^2 I)$, $w_t \sim \mathcal{N}(0, \sigma_w^2 I)$ for $t \geq 0$ as described in Section II. Assume, we have an estimate $\hat{G}$ of $G$ that is independent of these variables and we employ the $y_t$ estimator $\hat{y}_t = \hat{G}\bar{u}_t$. Then, $\mathbb{E}[\|y_t - \hat{y}_t\|_{\ell_2}^2] \leq$

$$\sigma_w^2 \|F\|_F^2 + \sigma_u^2 \|G - \hat{G}\|_F^2 + m\sigma_z^2 + \|CA^{T-1}\|^2 \text{tr}(\Gamma_\infty).$$

*Proof:* Following (II.5), the key observation is that for a fixed $t$, $\bar{u}_t, \bar{w}_t, z_t, e_t$ are all independent and the

associated errors are uncorrelated. Since $\bar{u}_t \sim \mathcal{N}(0, \sigma_u^2 I)$, $\mathbb{E}[\|(G - \hat{G})\bar{u}\|_{\ell_2}^2] = \sigma_u^2 \|G - \hat{G}\|_F^2$. Same argument applies to $\bar{w} \sim \mathcal{N}(0, \sigma_w^2 I), z_t \sim \mathcal{N}(0, \sigma_z^2 I)$ and $e_t$ which obeys $\mathbb{E}[\|e_t\|_{\ell_2}^2] = \text{tr}(\Sigma(e_t))$. Observe that $i$th largest eigenvalue $\lambda_i(\Sigma(e_t))$ of $\Sigma(e_t)$ is upper bounded by $\|CA^{T-1}\|^2 \lambda_i(\Sigma(x_{t-T+1}))$ via Min-Max principle [37] hence $\mathbb{E}[\|e_t\|_{\ell_2}^2] \leq \|CA^{T-1}\|^2 \text{tr}(\Sigma(x_{t-T+1})) \leq \|CA^{T-1}\|^2 \text{tr}(\Gamma_\infty)$. ∎

## IV. MARKOV PARAMETERS TO HANKEL MATRIX: LOW ORDER APPROXIMATION OF STABLE SYSTEMS

So far our attention has focused on estimating the impulse response $G$ for a particular horizon $T$. Clearly, we are also interested in understanding how well we learn the overall behavior of the system by learning a finite impulse approximation. In this section, we will apply our earlier results to approximate the overall system by using as few samples as possible. A useful idea towards this goal is taking advantage of the stability of the system. The Markov parameters decay exponentially fast if the system is stable i.e. $\rho(A) < 1$. This means that, most of the Markov parameters will be very small after a while and not learning them might not be a big loss for learning the overall behavior. In particular, $\tau$'th Markov parameter obeys

$$\|CA^\tau B\| \leq \Phi(A)\rho(A)^{\tau/2}\|C\|\|B\|.$$

This implies that, the impact of the impulse response terms we don't learn can be upper bounded. For instance, the total spectral norm of the tail terms obey

$$\sum_{\tau=T-1}^\infty \|CA^\tau B\| \leq \sum_{\tau=T-1}^\infty \Phi(A)\rho(A)^{\tau/2}\|C\|\|B\|$$
$$\leq \frac{\Phi(A)\|C\|\|B\|\rho(A)^{(T-1)/2}}{1 - \rho(A)^{1/2}}. \qquad \text{(IV.1)}$$

To proceed fix a finite horizon $K \geq T$ that will later be allowed to go infinity. Represent the estimate $\hat{G}$ as $[\hat{D}, \hat{G}_0, \ldots \hat{G}_{T-2}]$ where $\hat{G}_i$ corresponds to the noisy estimate of $CA^i B$. Now, let us consider the estimated and true order $K$ Markov parameters

$$\hat{G}^{(K)} = [\hat{D}, \hat{G}_0, \ldots \hat{G}_{T-2}\ 0\ \ldots\ 0]$$
$$G^{(K)} = [D, CB, CAB \ldots CA^{K-2}B].$$

We can construct Hankel matrices from these as follows.

*Definition 4.1 (Hankel matrix):* Given a block matrix $X = [X_1, X_2, \ldots X_T] \in \mathbb{R}^{m \times Tp}$ and integers $T_1, T_2$ satisfying $T_1 + T_2 \leq T$, define the associated $(T_1, T_2)$ block Hankel matrix $H = H(X) \in \mathbb{R}^{T_1 m \times T_2 p}$ to be the $T_1 \times T_2$ block matrix with $m \times p$ size blocks where $(i, j)$th block is given by

$$H[i, j] = X_{i+j}.$$

Observe that $H$ does not contain $X_1$, which shall correspond to the $D$ (or $\hat{D}$) matrix for our purposes. This is solely for notational convenience in the next section where the goal is identifying the $A, B, C$ matrices.

**Hankel operator:** Following Definition 4.1, we can create $K \times K$ block matrices $H^{(K)} = H(G^{(2K)})$ and $\hat{H}^{(K)} = H(\hat{G}^{(2K)})$ of size $mK \times pK$. The Hankel operator is the infinite dimensional linear operator obtained by $H^{(\infty)} = \lim_{K \to \infty} H^{(K)}$ and is critical for control applications. The

following theorem merges results of this section with a specific choice of $T$ to provide approximation bounds for the infinite Markov operator $\boldsymbol{G}^{(\infty)}$ and Hankel operator $\boldsymbol{H}^{(\infty)}$. For notational simplicity, we shall assume that there is no process noise.

*Theorem 4.2:* Suppose the spectral radius obeys $\rho(\boldsymbol{A}) < 1$. Fix a precision $1 > \varepsilon_0 > 0$ and suppose there is no process noise. Assume sample size $N$ and estimation horizon $T$ satisfy[3]

$$N / \log^2(N) \gtrsim cTp \log^2(Tp)$$
$$T \gtrsim \frac{c_0 + \log(N/T + T(1 + m/p)) - \log \varepsilon_0}{-\log \rho(\boldsymbol{A})}. \qquad \text{(IV.2)}$$

Then, given observations of a single trajectory until time $\bar{N} = N + T - 1$ and estimating first $T$ Markov parameters via least-squares estimator (II.7), with high probability, the following bounds hold on the infinite impulse response and Hankel matrix of the system:

$$\|\boldsymbol{G}^{(\infty)} - \hat{\boldsymbol{G}}^{(\infty)}\| \le (8 \frac{\sigma_z}{\sigma_u} + \varepsilon_0) \sqrt{\frac{Tp + m}{N}}$$

$$\|\boldsymbol{H}^{(\infty)} - \hat{\boldsymbol{H}}^{(\infty)}\| \le (8 \frac{\sigma_z}{\sigma_u} + \varepsilon_0) \sqrt{T \times \frac{Tp + m}{N}}.$$

In essence, the theorem above is a corollary of Theorem 3.2. However, it further simplifies the bounds and also provides approximation to systems overall behavior (e.g. infinite Hankel matrix). In particular, these bounds exploit stability of the system and allows us to treat the system as if it has a logarithmic order. Observe that (IV.2) only logarithmically depends on the critical problem variables such as precision $\varepsilon_0$ and spectral radius. In essence, the effective system order is dictated by the eigen-decay and equal to $T \sim \mathcal{O}(-\frac{1}{\log(\rho(\boldsymbol{A}))})$ hence stability allows us to treat the system as if it has a logarithmically small order. Ignoring logarithmic terms except $\rho(\boldsymbol{A})$, using $\varepsilon_0, \sigma_z/\sigma_u = \mathcal{O}(1)$ and picking

$$T = \mathcal{O}(\frac{-1}{\log(\rho(\boldsymbol{A}))}) \quad \text{and} \quad N = \mathcal{O}(\delta^{-2}(Tp + m)),$$

guarantees

$$\|\boldsymbol{G}^{(\infty)} - \hat{\boldsymbol{G}}^{(\infty)}\| \le \delta \quad \text{and} \quad \|\boldsymbol{H}^{(\infty)} - \hat{\boldsymbol{H}}^{(\infty)}\| \lesssim \frac{-\delta}{\log(\rho(\boldsymbol{A}))}.$$

By invoking standard results (see, e.g., [28]) on the singular values of the infinite Hankel matrices, one can also upper and lower bound the $\mathcal{H}_\infty$-norm error between the true system and the computed approximation. Observe that, this sample size bound is independent of the state dimension $n$ and only linearly grows with $p$ however it exhibits dependence on the spectral radius. It is useful in the regime when the system order is high and its Markov parameters decay rapidly. In scenarios where the system order is small and the spectral radius is rather large, choosing $T$ to be the true system order and finding a low-order realization, as discussed in the next section, would provide more informative bounds.

## V. NON-ASYMPTOTIC SYSTEM IDENTIFICATION VIA HO-KALMAN

In this section, we first describe the Ho-Kalman algorithm [26] that generates $\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}, \boldsymbol{D}$ from the Markov parameter

---

**Algorithm 1** Ho-Kalman Algorithm to find a State-Space Realization.

1: **procedure** HO-KALMAN MINIMUM REALIZATION
2: **Inputs:** Length $T$, Markov parameter matrix estimate $\hat{\boldsymbol{G}}$, system order $n$, Hankel shape $(T_1, T_2 + 1)$ with $T_1 + T_2 + 1 = T$, and $T_1, T_2 \ge n$.
3: **Outputs:** State-space realization $\hat{\boldsymbol{A}}, \hat{\boldsymbol{B}}, \hat{\boldsymbol{C}}$.
4:    Form the Hankel matrix $\hat{\boldsymbol{H}} \in \mathbb{R}^{mT_1 \times p(T_2+1)}$ from $\hat{\boldsymbol{G}}$.
5:    $\hat{\boldsymbol{H}}^- \in \mathbb{R}^{mT_1 \times pT_2} \leftarrow$ first-$pT_2$-columns-of-$(\hat{\boldsymbol{H}})$.
6:    $\hat{\boldsymbol{L}} \in \mathbb{R}^{mT_1 \times pT_2} \leftarrow$ rank-$n$-approximation-of-$(\hat{\boldsymbol{H}}^-)$.
7:    $\boldsymbol{U}, \boldsymbol{\Sigma}, \boldsymbol{V} = \text{SVD}(\hat{\boldsymbol{L}})$.
8:    $\hat{\boldsymbol{O}} \in \mathbb{R}^{mT_1 \times n} \leftarrow \boldsymbol{U}\boldsymbol{\Sigma}^{1/2}$.
9:    $\hat{\boldsymbol{Q}} \in \mathbb{R}^{n \times pT_2} \leftarrow \boldsymbol{\Sigma}^{1/2}\boldsymbol{V}^*$.
10:    $\hat{\boldsymbol{C}} \leftarrow$ first-$m$-rows-of-$(\hat{\boldsymbol{O}})$.
11:    $\hat{\boldsymbol{B}} \leftarrow$ first-$p$-columns-of-$(\hat{\boldsymbol{Q}})$.
12:    $\hat{\boldsymbol{H}}^+ \in \mathbb{R}^{mT_1 \times pT_2} \leftarrow$ last-$pT_2$-columns-of-$(\hat{\boldsymbol{H}})$.
13:    $\hat{\boldsymbol{A}} \leftarrow \hat{\boldsymbol{O}}^\dagger \hat{\boldsymbol{H}}^+ \hat{\boldsymbol{Q}}^\dagger$.
14: **return** $\hat{\boldsymbol{A}} \in \mathbb{R}^{n \times n}, \hat{\boldsymbol{B}} \in \mathbb{R}^{n \times p}, \hat{\boldsymbol{C}} \in \mathbb{R}^{m \times n}$.
15: **end procedure**

---

matrix $\boldsymbol{G}$. We also show that the algorithm is stable to perturbations in $\boldsymbol{G}$ and the output of Ho-Kalman gracefully degrades as a function of $\|\boldsymbol{G} - \hat{\boldsymbol{G}}\|$. Combining this with Theorem 3.1 implies *guaranteed* non-asymptotic identification of multi-input multi-output systems from a *single trajectory*. We remark that results of this section do not assume stability and applies to arbitrary, possibly unstable, systems.

### A. System Identification Algorithm

Given a noisy estimate $\hat{\boldsymbol{G}}$ of $\boldsymbol{G}$, we wish to learn good system matrices $\hat{\boldsymbol{A}}, \hat{\boldsymbol{B}}, \hat{\boldsymbol{C}}, \hat{\boldsymbol{D}}$ from $\hat{\boldsymbol{G}}$ up to trivial ambiguities. This will be achieved by using Algorithm 1 which admits the matrix $\hat{\boldsymbol{G}}$, system order $n$ and Hankel dimensions $T_1, T_2$ as inputs. Throughout this section, we make the following two assumptions to ensure that the system we wish to learn is order-$n$ and our system identification problem is well-conditioned.

- the system is observable and controllable; hence $n > 0$ is the order of the system.
- $(T_1, T_2)$ Hankel matrix $\boldsymbol{H}(\boldsymbol{G})$ formed from $\boldsymbol{G}$ is rank-$n$. This can be ensured by choosing sufficiently large $T_1, T_2$. In particular $T_1 \ge n, T_2 \ge n$ is guaranteed to work by the first assumption above.

Learning state-space representations is a non-trivial, inherently non-convex problem. Observe that there are multiple state-space realizations that yields the same system and Markov matrix $\boldsymbol{G}$. In particular, for any nonsingular matrix $\boldsymbol{T} \in \mathbb{R}^{n \times n}$,

$$\boldsymbol{A}' = \boldsymbol{T}^{-1}\boldsymbol{A}\boldsymbol{T}, \ \boldsymbol{B}' = \boldsymbol{T}^{-1}\boldsymbol{B}, \ \boldsymbol{C}' = \boldsymbol{C}\boldsymbol{T},$$

is a valid realization and yields the same system. Hence, similarity transformations of $\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}$ generate a class of solutions. Note that $\boldsymbol{D}$ is already estimated as part of $\boldsymbol{G}$. Since $\boldsymbol{D}$ is a submatrix of $\boldsymbol{G}$, we clearly have

$$\|\boldsymbol{D} - \hat{\boldsymbol{D}}\| \le \|\boldsymbol{G} - \hat{\boldsymbol{G}}\|.$$

Hence, we focus our attention on learning $\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}$. Suppose we have access to the true Markov parameters $\boldsymbol{G}$ and the

---

[3]Exact form of the bounds depend on $\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}$ and is provided in the proof.

This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication. Citation information: DOI 10.1109/TAC.2021.3083651, IEEE Transactions on Automatic Control

6     GENERIC COLORIZED JOURNAL, VOL. XX, NO. XX, XXXX 2017

corresponding $(T_1, T_2 + 1)$ Hankel matrix $H(G)$. In this case, $H$ is a rank-$n$ matrix and $(i,j)$th block of $H$ is equal to $CA^{i+j-2}B$. Defining (extended) controllability and observability matrices $Q = [B, \; AB, \; \dots \; A^{T_2}B]$ and $O = [C^*, \; (CA)^*, \; \dots \; (CA^{T_1-1})^*]^*$, we have $H = OQ$. However, it is not clear how to find $O, Q$.

The Ho-Kalman algorithm accomplishes this task by finding a balanced realization and returning *some* $\hat{A}, \hat{B}, \hat{C}$ matrices from possibly noisy Markov parameter matrix $\hat{G}$. Let the input to the algorithm be $\hat{G} = [\hat{D}, \; \hat{G}_0, \; \dots \; \hat{G}_{T-2}]$ where $\hat{G}_i$ corresponds to the noisy estimate of $CA^iB$. We construct the $(T_1, T_2+1)$ Hankel matrix $\hat{H}$ as described above so that $(i,j)$th block of $\hat{H}$ is equal to $\hat{G}_{i+j-2}$. Let $\hat{H}^- \in \mathbb{R}^{mT_1 \times pT_2}$ be the submatrix of $\hat{H}$ after discarding the rightmost $mT_1 \times p$ block and $\hat{L}$ be the best rank-$n$ approximation of $\hat{H}^-$ obtained by setting its all but top $n$ singular values to zero. Let $\hat{H}^+$ be the submatrix after discarding the left-most $mT_1 \times p$ block. Note that both $\hat{L}, \hat{H}^+$ have size $\mathbb{R}^{mT_1 \times pT_2}$. Take the singular value decomposition (SVD) of the rank-$n$ matrix $\hat{L}$ as $\hat{L} = U\Sigma V^*$ (with $\Sigma \in \mathbb{R}^{n \times n}$) and write

$$\hat{L} = (U\Sigma^{1/2})\Sigma^{1/2}V^* = \hat{O}\hat{Q}.$$

If $\hat{G}$ was equal to the ground truth $G$, then $\hat{O}, \hat{Q}$ would correspond to the order $T_1$ observability matrix $\bar{O} = U\Sigma^{1/2}$ and the order $T_2$ controllability matrix $\bar{Q} = \Sigma^{1/2}V^*$ of the actual balanced realization based on *noiseless* SVD. Here, $\bar{O}, \bar{Q}$ matrices are not necessarily equal to $O, Q$, however they yield the same system. Note that, the columns of $\hat{O}, \hat{Q}$ are the scaled versions of the left and right singular vectors of $\hat{L}$ respectively. The Ho-Kalman algorithm finds $\hat{A}, \hat{B}, \hat{C}$ as follows.

- $\hat{C}$ is the first $m \times n$ submatrix of $\hat{O}$.
- $\hat{B}$ is the first $n \times p$ submatrix of $\hat{Q}$.
- $\hat{A} = \hat{O}^\dagger \hat{H}^+ \hat{Q}^\dagger$.

This procedure (Ho-Kalman) returns the true balanced realization of the system when Markov parameters are known i.e. $\hat{G} = G$. Our goal is to show that even with noisy Markov parameters, this procedure returns good estimates of the true balanced realization. We remark that there are variations of this procedure; however the core idea is the same and they are equivalent when the true Markov parameters are used as input [2]. For instance, when constructing $\hat{H}$, one can attempt to improve the noise robustness of the algorithm by picking balanced dimensions $mT_1 \approx pT_2$.

### B. Robustness of the Ho-Kalman Algorithm

Observe that $\hat{H}, \hat{H}^-, \hat{L}, \hat{H}^+, \hat{O}, \hat{Q}$ of Algorithm 1 are functions of the input matrix $\hat{G}$. For the subsequent discussion, we let

- $H, H^-, L, H^+, O, Q$ be the matrices corresponding to ground truth $G$.
- $\hat{H}, \hat{H}^-, \hat{L}, \hat{H}^+, \hat{O}, \hat{Q}$ be the matrices corresponding to the estimate $\hat{G}$.

Furthermore, let $\bar{A}, \bar{B}, \bar{C}$ be the actual balanced realization associated with $G$ and let $\hat{A}, \hat{B}, \hat{C}$ be the Ho-Kalman output associated with $\hat{G}$. Note that $L = H^-$ since $H^-$ is already rank $n$. We now provide a lemma relating the estimation error of $G$ to that of $L$ and $H$.

*Lemma 5.1:* $H, \hat{H}$ and $L, \hat{L}$ satisfies the following perturbation bounds,

$$\max\{\|H^+ - \hat{H}^+\|, \|H^- - \hat{H}^-\|\}$$
$$\leq \|H - \hat{H}\| \leq \sqrt{\min\{T_1, T_2 + 1\}}\|G - \hat{G}\|. \quad \text{(V.1)}$$
$$\|L - \hat{L}\| \leq 2\|H^- - \hat{H}^-\| \leq 2\sqrt{\min\{T_1, T_2\}}\|G - \hat{G}\|. \quad \text{(V.2)}$$

Let us denote the $n$th largest singular value of $L$ via $\sigma_n(L)$. Note that $\sigma_n(L)$ is the smallest nonzero singular value of $L$ since rank$(L) = n$. A useful implication of Theorem 3.1 (in light of Lemma 5.1) is that if $\sigma_n(L)$ is large enough, the true system order $n$ can be non-asymptotically estimated from the noisy Markov parameter estimates via singular value thresholding.

Our next result shows the robustness of the Ho-Kalman algorithm to possibly adversarial perturbations on the Markov parameter matrix $G$.

*Theorem 5.2:* Suppose $H$ and $\hat{H}$ be the Hankel matrices derived from $G$ and $\hat{G}$ respectively per Definition 4.1. Let $\bar{A}, \bar{B}, \bar{C}$ be the state-space realization corresponding to the output of Ho-Kalman with input $G$ and $\hat{A}, \hat{B}, \hat{C}$ be the state-space realization corresponding to output of Ho-Kalman with input $\hat{G}$. Suppose the system $A, B, C, D$ is *observable and controllable* and let $O, Q$ and $\hat{O}, \hat{Q}$ be order-$n$ controllability/observability matrices associated with $G$ and $\hat{G}$ respectively. Suppose $\sigma_n(L) > 0$ and perturbation obeys

$$\|L - \hat{L}\| \leq \sigma_n(L)/2. \quad \text{(V.3)}$$

Then, there exists a unitary matrix $T \in \mathbb{R}^{n \times n}$ such that,

$$\|\bar{C} - \hat{C}T\|_F^2 \leq \|O - \hat{O}T\|_F^2 \leq 10n\|L - \hat{L}\|^2/\sigma_n(L), \quad \text{(V.4)}$$
$$\|\bar{B} - T^*\hat{B}\|_F^2 \leq \|Q - T^*\hat{Q}\|_F^2 \leq 10n\|L - \hat{L}\|^2/\sigma_n(L). $$

Furthermore, hidden state matrices $\hat{A}, \bar{A}$ satisfy

$$\|\bar{A} - T^*\hat{A}T\|_F \leq \frac{9\sqrt{n}}{\sigma_n(L)}\left(\frac{\|L - \hat{L}\|}{\sigma_n(L)}\|H^+\| + \|H^+ - \hat{H}^+\|\right). \quad \text{(V.5)}$$

Above, $\|H^+ - \hat{H}^+\|, \|L - \hat{L}\|$ are perturbation terms that can be bounded in terms of $\|H - \hat{H}\|$ or $\|G - \hat{G}\|$ via Lemma 5.1. This result shows that Ho-Kalman solution is robust to noise up to trivial ambiguities. Robustness is controlled by $\sigma_n(L)$ which corresponds to the weakest mode of the system. This is not surprising since "weakest" here is in terms of controllability and observability, therefore $\sigma_n(L)$ being small indicates that there is a mode of the system that is hard to identify. We remark that, for a stable system and for reasonably large $T_2$ choice, we have that $\sigma_n(L) \approx \sigma_n(H)$. This is because $L = H^-$ is obtained by discarding the last block column of $H$ which is exponentially small in $T_2$. Additionally, observe that, $\sigma_n(L)$ implicitly depends on the Hankel size $(T_1, T_2)$. However, as $(T_1, T_2)$ grows larger, it will similarly quickly converge to the system dependent quantity $\sigma_n(H^{(\infty)})$ where $H^{(\infty)}$ is the infinite Hankel matrix of the system which was introduced in Section IV. Finally, we note that, this bound is consistent with that of [32] and would lead to the ideal finite sample estimation error rate of $1/\sqrt{N}$ for the state-space matrices.

Since the Ho-Kalman algorithm is based on SVD, having a good control over singular vectors is crucial for the proof.

We do this by utilizing the perturbation results from the recent literature [38]. While we believe our result has the correct dependency, it is in terms of Frobenius norm rather than spectral. Having a better spectral norm control over $\bar{A}, \bar{B}, \bar{C}$ would be an ideal future improvement.

Our next result combines Ho-Kalman's robustness with finite sample learning bounds of Theorem 3.1 to establish an end-to-end estimation guarantee on the state-space matrices.

*Theorem 5.3:* Consider the setups of Theorems 3.1 and 5.2. Suppose $\sigma_u = 1$ and $\sigma_z, \sigma_e, \sigma_w \|F\|$ are bounded above by constants. Set $N_0 = Tq \log^2(Tq)$ where $q = p + m + n$. Suppose $\sigma_n(L) > 0$ and sample size obeys

$$N / \log^2(Nq) \gtrsim TN_0 / \sigma_n^2(L). \tag{V.6}$$

Then, with high probability (same as Thm 3.1), there exists a unitary matrix $T \in \mathbb{R}^{n \times n}$ such that,

$$\max\{\|\bar{C} - \hat{C}T\|_F, \|O - \hat{O}T\|_F, \|\bar{B} - T^*\hat{B}\|_F, \|Q - T^*\hat{Q}\|_F\}$$

$$\leq \frac{\sqrt{CnT} \log(Nq)}{\sqrt{\sigma_n(L)}} \sqrt{\frac{N_0}{N}}. \tag{V.7}$$

Furthermore, hidden state matrices $\hat{A}, \bar{A}$ satisfy

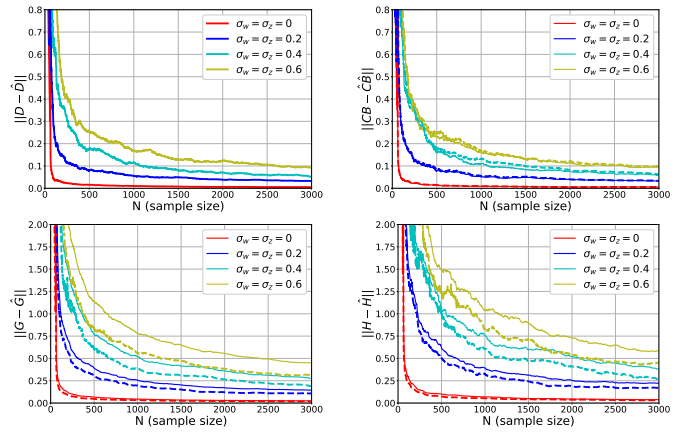$$\|\bar{A} - T^*\hat{A}T\|_F \leq \frac{C\sqrt{nT} \log(Nq)\|H\|}{\sigma_n^2(L)} \sqrt{\frac{N_0}{N}}. \tag{V.8}$$

This theorem achieves a standard $1/\sqrt{N}$ estimation rate on the state-space. Recall that, ignoring log factors, $N_0 \sim Tq$. This implies that, in terms of problem dimensions, we essentially need $N \gtrsim T^2 qn$ samples (to achieve constant error). We suspect that quadratic dependence on $T$ is a proof artifact which can hopefully be refined in future works.

## VI. NUMERICAL EXPERIMENTS

We consider MIMO (multiple input, multiple output) systems with $m = 2$ sensors, $n = 5$ hidden states and input dimension $p = 3$. To assess the typical performance of the least-squares and the Ho-Kalman algorithms, we consider a random state-space models as follows. For each realization of a system trajectory, we generate different $A, B, C, D$ matrices which are drawn with i.i.d. normally distributed entries. The entries of $C, D$ have variance $1/m$ and $B$ has variance $1/n$ to ensure these matrices are isometric in the sense that they approximately preserve norm of the input vector. Hence, the impact of the standard deviations $\sigma_u, \sigma_w, \sigma_z$ are properly normalized. The input variance is fixed at $\sigma_u = 1$ however noise variances will be modified during the experiments.

The most critical component of an LTI system is the $A$ matrix. We generate $A$'s entries with variance 1. We then take an eigenvalue decomposition and scale eigenvalues of $A$ to have absolute value $\rho = 0.85$. The upper bound $\rho$ implies that we are working with stable matrices and the effect of unknown state vanishes for large $T$. Eventually $A$ have four complex and one real eigenvalue and three out of five eigenvalues are equal to $\rho$ in absolute value.

In our experiments we set $T = 20$ and work with a $10 \times 10$ block Hankel in Algorithm 1. This means that $G$ is of size $2 \times 60$ and $H$ is of size $20 \times 30$. We pick noise configurations $\sigma_w, \sigma_z$ and generate a single rollout of the system until time



Fig. 1: We consider the matrices that can directly be inferred from the Markov parameter matrix $G$. These are $D, CB$ which are the first two block submatrices of $G$, $G$ itself, and $H$ which is the Hankel matrix that is constructed from blocks of $G$. These results are for $T = 20$ which implies $G \in \mathbb{R}^{2 \times 60}$ and $H \in \mathbb{R}^{20 \times 30}$ as we picked $T_1 = T_2 + 1 = 10$. The solid lines are the estimates obtained from $\hat{G}$ whereas the dashed lines are the estimates obtained from the Ho-Kalman output $\hat{A}, \hat{B}, \hat{C}$. Ho-Kalman leads to noticeable improvement over simply using $\hat{G}$ when estimating the Hankel matrix.

$\bar{N}_{\max} = 3000$. For each $\bar{N} \leq \bar{N}_{\max}$, we solve the system via (II.7) to obtain the $\hat{G}$ and use Algorithm 1 to obtain a state-space realization $\hat{A}, \hat{B}, \hat{C}, \hat{D}$. The $x$-axis displays $N$ (the amount of available data at time $t = \bar{N}$) and the $y$-axis displays the estimation error. Each curve in the figures is generated by averaging the outcomes of 20 independent realizations of single trajectories each generated from a random $(A, B, C, D)$ choice described above.

Importantly, at the rare event that the estimated $\hat{A}$ has an eigenvalue larger than $0.95$ (in absolute value), we scale them to be $0.95$. This operation enforces the stability of the system estimate. While we verified that $\rho(\hat{A}) \geq 1$ rarely happens for large $N$, clipping smooths out the results by discarding outliers due to the exponents of $\hat{A}$ appearing in $G$ and $H$.

In Figure 1, we investigate the problem of estimating the matrices $D, CB, G, H$. $D, CB$ are the first two impulse responses. Estimating $G$ and the associated Hankel matrix $H$ helps verify our findings in Theorem 3.2. The solid lines are the estimates obtain directly from Markov parameters $\hat{G}$. The dashed lines are the estimates obtained *after the Ho-Kalman* procedure i.e. constructed from $(\hat{A}, \hat{B}, \hat{C}, \hat{D})$.

We plotted curves for varying noise levels $\sigma_w = \sigma_z \in \{0, 0.2, 0.4, 0.6\}$. *The first major conclusion is that indeed estimation accuracy drastically improves as we observe the system for a longer period of time* and collect more data. Note that $D$ and $CB$ are submatrices of $G$ hence their associated spectral norm errors are strictly lower compared to $\|G - \hat{G}\|$. Per Definition 4.1, $H$ is constructed from the blocks of $G$ and its spectral norm error is in lines with $G$. *The second major conclusion is that Ho-Kalman procedure is indeed robust.* The dashed lines are in fact under the solid lines indicating that Ho-Kalman outputs a more refined system (compared to $\hat{G}$) by projecting initial Markov parameters to the set of low-order systems. Another observation is that estimation error grows gracefully as a function of the noise levels for all matrices of

interest. Since the chosen $T$ is large and $\rho = 0.85$ is reasonably stable, the error due to unknown initial conditions (i.e. $e_t$) is fairly negligible. Hence when $\sigma_w = \sigma_z = 0$, we quickly achieve near 0 estimation error as the impact of the $e_t$ term is small.

## VII. CONCLUSIONS

In this paper, we analyzed the sample complexity of linear system identification from input/output data. Our analysis neither requires multiple independent trajectories nor relies on splitting the trajectory into non-overlapping intervals, therefore makes very efficient use of the available data from a single trajectory. More crucially, it does not rely on state measurements and works with only the inputs and outputs. Based on this analysis, we showed that one can approximate system's Hankel operator using near optimal amount of samples and shed light on the robustness of finding a balanced realization. This type of analysis is particularly useful in understanding how certain system properties affect the learning rates. There are many directions for future work. We are especially interested in what type of recovery guarantees can be obtained if additional structural constraints, such as subspace constraints, on the system matrices are known [39].

## ACKNOWLEDGEMENTS

## REFERENCES

[1] L. Ljung, "System identification," in *Signal analysis and prediction*. Springer, 1998, pp. 163–173.

[2] D. Bauer and M. Jansson, "Analysis of the asymptotic properties of the moesp type of subspace algorithms," *Automatica*, vol. 36, no. 4, pp. 497–509, 2000.

[3] P. Van Overschee and B. De Moor, *Subspace identification for linear systems: Theory–Implementation–Applications*. Springer Science & Business Media, 2012.

[4] E. Weyer, R. C. Williamson, and I. M. Mareels, "Finite sample properties of linear model identification," *IEEE Transactions on Automatic Control*, vol. 44, no. 7, pp. 1370–1383, 1999.

[5] M. C. Campi and E. Weyer, "Finite sample properties of system identification methods," *IEEE Transactions on Automatic Control*, vol. 47, no. 8, pp. 1329–1334, 2002.

[6] H. Akçay, "The size of the membership-set in a probabilistic framework," *Automatica*, vol. 40, no. 2, pp. 253–260, 2004.

[7] M. Vidyasagar and R. Karandikar, "System identification: a learning theory approach," in *Decision and Control, 2001. Proceedings of the 40th IEEE Conference on*, vol. 2. IEEE, 2001.

[8] A. Carè, B. C. Csáji, M. C. Campi, and E. Weyer, "Finite-sample system identification: An overview and a new correlation method," *IEEE Control Systems Letters*, vol. 2, no. 1, pp. 61–66, 2018.

[9] M. Simchowitz, H. Mania, S. Tu, M. I. Jordan, and B. Recht, "Learning without mixing: Towards a sharp analysis of linear system identification," *arXiv preprint arXiv:1802.08334*, 2018.

[10] M. K. S. Faradonbeh, A. Tewari, and G. Michailidis, "Finite time analysis of optimal adaptive policies for linear-quadratic systems," *arXiv preprint arXiv:1711.07230*, 2017.

[11] E. Hazan, K. Singh, and C. Zhang, "Learning linear dynamical systems via spectral filtering," in *Advances in Neural Information Processing Systems*, 2017, pp. 6705–6715.

[12] M. Fazel, R. Ge, S. M. Kakade, and M. Mesbahi, "Global convergence of policy gradient methods for linearized control problems," *arXiv preprint arXiv:1801.05039*, 2018.

[13] J. Pereira, M. Ibrahimi, and A. Montanari, "Learning networks of stochastic differential equations," in *Advances in Neural Information Processing Systems*, 2010, pp. 172–180.

[14] M. Hardt, T. Ma, and B. Recht, "Gradient descent learns linear dynamical systems," *arXiv preprint arXiv:1609.05191*, 2016.

[15] S. Dean, H. Mania, N. Matni, B. Recht, and S. Tu, "On the sample complexity of the linear quadratic regulator," *arXiv preprint arXiv:1710.01688*, 2017.

[16] R. Boczar, N. Matni, and B. Recht, "Finite-data performance guarantees for the output-feedback control of an unknown system," *arXiv preprint arXiv:1803.09186*, 2018.

[17] S. Arora, E. Hazan, H. Lee, K. Singh, C. Zhang, and Y. Zhang, "Towards provable control for unknown linear dynamical systems," *ICLR workshop*, 2018.

[18] S. Tu, R. Boczar, A. Packard, and B. Recht, "Non-asymptotic analysis of robust control from coarse-grained identification," *arXiv preprint arXiv:1707.04791*, 2017.

[19] T. Sarkar and A. Rakhlin, "How fast can linear dynamical systems be learned?" *arXiv preprint arXiv:1812.01251*, 2018.

[20] M. K. S. Faradonbeh, A. Tewari, and G. Michailidis, "Finite time identification in unstable linear systems," *Automatica*, vol. 96, pp. 342–353, 2018.

[21] S. L. Tu, "Sample complexity bounds for the linear quadratic regulator," Ph.D. dissertation, UC Berkeley, 2019.

[22] M. S. Fledderjohn, M. S. Holzel, H. J. Palanthandalam-Madapusi, R. J. Fuentes, and D. S. Bernstein, "A comparison of least squares algorithms for estimating markov parameters," in *American Control Conference (ACC), 2010*. IEEE, 2010, pp. 3735–3740.

[23] R. E. Skelton and G. Shi, "The data-based LQG control problem," in *Decision and Control, 1994., Proceedings of the 33rd IEEE Conference on*, vol. 2. IEEE, 1994, pp. 1447–1452.

[24] K. Furuta and M. Wongsaisuwan, "Discrete-time LQG dynamic controller design using plant markov parameters," *Automatica*, vol. 31, no. 9, pp. 1317–1324, 1995.

[25] M. A. Santillo and D. S. Bernstein, "Adaptive control based on retrospective cost optimization," *Journal of guidance, control, and dynamics*, vol. 33, no. 2, pp. 289–304, 2010.

[26] B. Ho and R. E. Kálmán, "Effective construction of linear state-variable models from input/output functions," *at-Automatisierungstechnik*, vol. 14, no. 1-12, pp. 545–548, 1966.

[27] K. Glover, "All optimal hankel-norm approximations of linear multi-variable systems and their l,∞-error bounds," *International journal of control*, vol. 39, no. 6, pp. 1115–1193, 1984.

[28] R. S. Sanchez-Pena and M. Sznaier, *Robust systems theory and applications*. Wiley-Interscience, 1998.

[29] A. Chiuso and G. Picci, "The asymptotic variance of subspace estimates," *Journal of Econometrics*, vol. 118, no. 1-2, pp. 257–291, 2004.

[30] S. Oymak and N. Ozay, "Non-asymptotic identification of LTI systems from a single trajectory," in *2019 American Control Conference (ACC)*. IEEE, 2019, pp. 5655–5661.

[31] T. Sarkar, A. Rakhlin, and M. A. Dahleh, "Finite-time system identification for partially observed lti systems of unknown order," *arXiv preprint arXiv:1902.01848*, 2019.

[32] M. Simchowitz, R. Boczar, and B. Recht, "Learning linear dynamical systems with semi-parametric least squares," *arXiv preprint arXiv:1902.00768*, 2019.

[33] H. Mania, S. Tu, and B. Recht, "Certainty equivalent control of lqr is efficient," *arXiv preprint arXiv:1902.07826*, 2019.

[34] A. Tsiamis and G. J. Pappas, "Finite sample analysis of stochastic system identification," *arXiv preprint arXiv:1903.09122*, 2019.

[35] D. C. Montgomery, E. A. Peck, and G. G. Vining, *Introduction to linear regression analysis*. John Wiley & Sons, 2012, vol. 821.

[36] F. Krahmer, S. Mendelson, and H. Rauhut, "Suprema of chaos processes and the restricted isometry property," *Communications on Pure and Applied Mathematics*, vol. 67, no. 11, pp. 1877–1904, 2014.

[37] R. A. Horn, R. A. Horn, and C. R. Johnson, *Matrix analysis*. Cambridge university press, 1990.

[38] S. Tu, R. Boczar, M. Simchowitz, M. Soltanolkotabi, and B. Recht, "Low-rank solutions of linear matrix equations via procrustes flow," *arXiv preprint arXiv:1507.03566*, 2015.

This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication. Citation information: DOI 10.1109/TAC.2021.3083651, IEEE Transactions on Automatic Control

AUTHOR *et al.*: PREPARATION OF PAPERS FOR IEEE TRANSACTIONS AND JOURNALS (FEBRUARY 2017)
9

[39] S. Fattahi and S. Sojoudi, "Sample complexity of sparse system identification problem," *arXiv preprint arXiv:1803.07753*, 2018.

[40] F. M. Dopico, "A note on sin $\theta$ theorems for singular subspace variations," *BIT Numerical Mathematics*, vol. 40, no. 2, pp. 395–403, 2000.

[41] Y. Yu, T. Wang, and R. J. Samworth, "A useful variant of the davis–kahan theorem for statisticians," *Biometrika*, vol. 102, no. 2, pp. 315–323, 2014.

[42] L. Meng and B. Zheng, "The optimal perturbation bounds of the moore–penrose inverse under the frobenius norm," *Linear Algebra and its Applications*, vol. 432, no. 4, pp. 956–963, 2010.

[43] P.-Å. Wedin, "Perturbation theory for pseudo-inverses," *BIT Numerical Mathematics*, vol. 13, no. 2, pp. 217–232, 1973.

[44] M. Talagrand, *The generic chaining: upper and lower bounds of stochastic processes*. Springer Science & Business Media, 2006.

[45] R. Adamczak *et al.*, "A note on the Hanson-Wright inequality for random vectors with dependencies," *Electronic Communications in Probability*, vol. 20, 2015.

[46] M. Rudelson, R. Vershynin *et al.*, "Hanson-wright inequality and sub-gaussian concentration," *Electronic Communications in Probability*, vol. 18, 2013.

[47] L. Hogben, *Handbook of linear algebra*. CRC Press, 2006.

## APPENDIX I
## PROOF OF THE RESULTS ON LEARNING MARKOV PARAMETERS

Our basic proof idea is to bound the individual terms that appear in the bound on the spectral norm of the error in the Markov parameter estimates given in Eq. (III.5). We first prove Theorem 3.2 which is our main theorem. It will be followed by the proof of Theorem 3.1.

### A. Proof of Theorem 3.2

*Proof:* The proof is obtained by combining estimates from the subsequent sections. Set $\Theta = \log^2(Tp) \log^2(Np)$ to simplify the notation. Picking $c \geq (\log 2)^{-4}$, our assumption of $N \geq cTp\Theta$ implies $N \geq T$ and $N \geq (\bar{N}+1)/2$ (using $\bar{N} = N + T - 1$). Consequently, $\log^2(\bar{N}p) \leq 4\log^2(Np)$ and we have $N \geq (c/4)Tp \log^2(Tp) \log^2(\bar{N}p)$. This fact will be useful when we need to utilize results of Appendix III. We first address the $\boldsymbol{Z}$ component of the error which is rather trivial to bound.

*Lemma 1.1:* Let $\boldsymbol{M} \in \mathbb{R}^{m \times n}$ be a tall matrix ($m \geq n$) with $\|\boldsymbol{M}\| \leq \eta$. Let $\boldsymbol{G} \in \mathbb{R}^{m \times k}$ be a matrix with independent standard normal entries. Then, with probability at least $1 - 2\exp(-t^2/2)$,

$$\|\boldsymbol{M}^*\boldsymbol{G}\| \leq \eta(\sqrt{2(n+k)} + t).$$

In particular, setting $t = \sqrt{2(n+k)}$, we find $\|\boldsymbol{M}^*\boldsymbol{G}\| \leq \eta\sqrt{8(n+k)}$ with probability at least $1 - 2\exp(-(n+k))$.

*Proof:* Suppose $\boldsymbol{M}$ have singular value decomposition $\boldsymbol{M} = \boldsymbol{V}_1 \boldsymbol{\Sigma} \boldsymbol{V}_2^*$ where $\boldsymbol{V}_1 \in \mathbb{R}^{m \times n}$. Observe that $\bar{\boldsymbol{G}} = \boldsymbol{V}_1^* \boldsymbol{G} \in \mathbb{R}^{n \times k}$ have i.i.d. $\mathcal{N}(0,1)$ entries. Also $\mathbb{E}[\|\bar{\boldsymbol{G}}\|] \leq \sqrt{n} + \sqrt{k} \leq \sqrt{2(n+k)}$. Applying Lipschitz Gaussian concentration on spectral norm, with probability at least $1 - 2\exp(-t^2/2)$, we obtain the relations $\|\boldsymbol{M}^*\boldsymbol{G}\| = \|\boldsymbol{V}_2 \boldsymbol{\Sigma} \bar{\boldsymbol{G}}\| = \|\boldsymbol{\Sigma}\bar{\boldsymbol{G}}\| \leq \eta(\sqrt{2(n+k)} + t)$ concluding the proof. ∎

The following corollary states the estimation error due to measurement noise ($\boldsymbol{Z}$ term).

*Corollary 1.2:* Let $\boldsymbol{U} \in \mathbb{R}^{N \times Tp}$ be the data matrix as in (II.6) and let $\boldsymbol{Z} \in \mathbb{R}^{N \times m}$ be the measurement noise matrix from (III.3). Suppose $N \geq cTp\Theta$ for some absolute constant $c > 0$. With probability at least $1 - 2\exp(-(Tp+m)) - \exp(-\Theta)$,

$$\|\boldsymbol{U}^*\boldsymbol{Z}\| \leq 4\sigma_u \sigma_z \sqrt{N(Tp+m)}.$$

*Proof:* Set $\eta = \sqrt{2N}\sigma_u$. Using $\bar{N} \geq N$, Lemma 3.3 yields

$$\mathbb{P}(\|\boldsymbol{U}\| \leq \eta) \geq 1 - \exp(-\Theta). \qquad (\text{I.1})$$

Hence, combining Lemmas 3.3 and 1.1, using the fact that $\boldsymbol{Z}, \boldsymbol{U}$ are independent, and adjusting for $\boldsymbol{Z}$'s variance $\sigma_z$, we find the result. ∎

Next, we apply Lemmas 3.3 and 3.4 to find that, for sufficiently large $c > 0$, whenever $N \geq cTp\Theta$,

$$\|(\boldsymbol{U}^*\boldsymbol{U})^{-1}\| \leq 2\sigma_u^{-2}/N, \qquad (\text{I.2})$$

$$\|\boldsymbol{U}^*\boldsymbol{W}\| \leq \frac{1}{2}\sigma_u\sigma_w \max\{\sqrt{N_w N}, N_w\},$$

where $N_w = cTq \log^2(Tq)\log^2(Nq)$ and $q = p + n$ with probability at least $1 - 2\exp(-\Theta)$. Finally, applying Theorem 4.1 with $\gamma = \frac{\|\boldsymbol{\Gamma}_\infty\|\Phi(\boldsymbol{A})^2\|\boldsymbol{C}\boldsymbol{A}^{T-1}\|^2}{1-\rho(\boldsymbol{A})^T}$ as in (IV.1), with probability at least $1 - T(\exp(-100Tp) + 2\exp(-100m))$,

$$\|\boldsymbol{U}^*\boldsymbol{E}\| \leq c_3\sigma_u\sqrt{T\max\{N, \frac{mT}{1-\rho(\boldsymbol{A})^{T/2}}\}\max\{Tp,m\}\gamma}.$$

Combining all of the estimates above via union bound and substituting $\theta$, with probability at least,

$$1 - 2\exp(-(Tp+m)) - 3(Np)^{-\log(Np)\log^2(Tp)}$$
$$- T(\exp(-100Tp) + 2\exp(-100m)),$$

the error term $\|\boldsymbol{G} - \hat{\boldsymbol{G}}\|$ of (III.5) is upper bounded by $\frac{R_z + R_e + R_w}{\sigma_u\sqrt{N}}$ where

$$R_z = 8\sigma_z\sqrt{Tp+m}, \qquad (\text{I.3})$$

$$R_w = \sigma_w\|\boldsymbol{F}\|\max\{\sqrt{N_w}, N_w/\sqrt{N}\}, \qquad (\text{I.4})$$

$$R_e = 2C\sqrt{(1 + \frac{mT}{N(1-\rho(\boldsymbol{A})^{T/2})})(Tp+m)T\gamma}. \qquad (\text{I.5})$$

Absorbing the $\times 2$ multiplier of $R_e$ into $C$ and observing $T\gamma = \sigma_e^2$, we conclude with the desired result. ∎

### B. Proof of Theorem 3.1

The proof uses the same strategy in Section I-A with slight modifications. We will repeat the argument for the sake of completeness. First of all, we utilize the same estimates based on Lemmas 3.3 and 3.4, namely (I.2) and (I.1) ($\|\boldsymbol{U}\| \leq \sqrt{2N}\sigma_u$) which hold with probability at least $1 - 3(Np)^{-\log(Np)\log^2(Tp)}$. Set $q' = p + n$. Since $q \geq q'$, observe that $N \geq N_0 \log^2(Nq) \geq N_w = cTq'\log^2(Tq')\log^2(Nq')$. Thus, we have that

$$\|\boldsymbol{U}^*\boldsymbol{W}\| \leq \frac{1}{2}\sigma_u\sigma_w\sqrt{N_w N} \leq \frac{1}{2}\sigma_u\sigma_w\sqrt{N_0 N}\log(Nq).$$

We use Lemma 1.1 with $t = \sqrt{2Tq}$ to obtain $\mathbb{P}(\|\boldsymbol{U}^*\boldsymbol{Z}\| \leq 4\sigma_u\sigma_z\sqrt{TqN}) \geq 1 - 2\exp(-Tq)$.

Finally, to bound the contribution of $\boldsymbol{E}$ we again apply Theorem 4.1. Since $\rho(\boldsymbol{A})^T \leq 0.99$, picking sufficiently large $c$, we observe that

$$\max\{N, \frac{mT}{1-\rho(\boldsymbol{A})^{T/2}}\} = N$$

Hence, setting $\sigma_e = \sqrt{\gamma T}$ with $\gamma$ defined in (IV.1) and applying Theorem 4.1 yields that for some $C > 0$

$$\|\boldsymbol{U}^*\boldsymbol{E}\| \leq C\sigma_u\sqrt{TN(Tp+m)\gamma} \leq C\sigma_u\sigma_e\sqrt{NTq},$$

holds with probability at least $1 - T(\exp(-100Tp) + 2\exp(-100m))$. Union bounding over all these events and following (III.5), with probability at least,

$$1 - 2\exp(-Tq) - 3(Np)^{-\log(Np)\log^2(Tp)}$$
$$- T(\exp(-100Tq) + 2\exp(-100m)),$$

we find the spectral norm estimation error of

$$\|\hat{\boldsymbol{G}} - \boldsymbol{G}\| \leq \frac{\frac{1}{2}\sigma_u\sigma_w\|\boldsymbol{F}\|\sqrt{N_0 N}\log(Nq)}{(\sigma_u^2 N)/2}$$
$$+ \frac{4\sigma_u\sigma_z\sqrt{TqN} + C\sigma_u\sigma_e\sqrt{TqN}}{(\sigma_u^2 N)/2}$$
$$\leq \frac{\sigma_w\|\boldsymbol{F}\|\sqrt{N_0}\log(Nq) + 8\sigma_z\sqrt{Tq} + 2C\sigma_e\sqrt{Tq}}{\sigma_u\sqrt{N}},$$

which is the desired bound after ensuring $\max\{8, 2C\}^2 Tq \leq N_0$ by picking $c$ (multiplier within $N_0$) to be sufficiently large.

### C. Proof of Theorem 4.2

*Proof:* To prove our bound, we will pick $T$ to be $T \geq \max(T_0, T_1, T_2, \frac{-2}{\log(\rho(\boldsymbol{A}))})$ for proper choices of $T_i$'s individually satisfying (IV.2). Let us start with $\boldsymbol{G}^{(\infty)}$ estimate. First observe that, the tail of the Markov parameters is bounded via (IV.1). Picking $T \geq T_1 := 1 - 2\frac{\log(2\varepsilon_0^{-1}(1-\rho(\boldsymbol{A})^{1/2})^{-1}\Phi(\boldsymbol{A})\|\boldsymbol{C}\|\|\boldsymbol{B}\|\sqrt{\frac{N}{Tp+m}})}{\log(\rho(\boldsymbol{A}))}$ implies that the right hand side of (IV.1) can be upper bounded as

$$\frac{\Phi(\boldsymbol{A})\|\boldsymbol{C}\|\|\boldsymbol{B}\|\rho(\boldsymbol{A})^{(T-1)/2}}{1 - \rho(\boldsymbol{A})^{1/2}} \leq \frac{1}{2}\varepsilon_0\sqrt{\frac{Tp+m}{N}} \iff$$
$$\rho(\boldsymbol{A})^{-(T-1)/2} \geq \frac{\Phi(\boldsymbol{A})\|\boldsymbol{C}\|\|\boldsymbol{B}\|\sqrt{\frac{N}{Tp+m}}}{\varepsilon_0(1 - \rho(\boldsymbol{A})^{1/2})} \quad \text{(I.6)}$$

Next, we will bound the spectral difference of order $T$ finite responses $\boldsymbol{G}$ and $\hat{\boldsymbol{G}}$. Let $T \geq -\frac{2}{\log(\rho(\boldsymbol{A}))}$ to ensure $\rho(\boldsymbol{A})^{T/2} \leq 1/2$. Applying Theorem 3.2, we will show that individual error summands due to $R_w, R_e, R_z$ are upper bounded. First, Theorem 3.2 is applicable due to the choice of $N$. $R_w$ summand is zero as $\sigma_w = 0$. Second, in order to bound the $R_e$ term, observe that, for some $C > 0$

$$\frac{R_e}{\sigma_u\sqrt{N}} \leq \frac{C}{4}\sigma_e\sqrt{\left(1 + \frac{mT}{N(1-\rho(\boldsymbol{A})^{T/2})}\right)(Tp+m)}$$
$$\leq \frac{\sigma_e}{\sigma_u}\frac{C}{4}\sqrt{1 + \frac{m}{p}}\sqrt{\frac{Tp+m}{N}}.$$

where we used $\sqrt{1 + 2mT/N} \leq \sqrt{1 + \frac{m}{p}}$. Since $\sigma_w = 0$, define,

$$\bar{\boldsymbol{\Gamma}}_\infty = \frac{\boldsymbol{\Gamma}_\infty}{\sigma_u^2} = \sum_{i=0}^{\infty} \boldsymbol{A}^i \boldsymbol{B}\boldsymbol{B}^*(\boldsymbol{A}^*)^i.$$

Defining $T_2 := 1 - 2\frac{\log(C\varepsilon_0^{-1}\Phi(\boldsymbol{A})^2\|\boldsymbol{C}\|\sqrt{T\|\bar{\boldsymbol{\Gamma}}_\infty\|(1+m/p)})}{\log(\rho(\boldsymbol{A}))}$ and setting $T \geq T_2$ guarantees the following upper bound on $\sigma_e$

$$\frac{\sigma_e}{\sigma_u} \leq 2\Phi(\boldsymbol{A})\|\boldsymbol{C}\boldsymbol{A}^{T-1}\|\sqrt{T\|\bar{\boldsymbol{\Gamma}}_\infty\|}$$
$$\leq 2\Phi(\boldsymbol{A})^2\|\boldsymbol{C}\|\rho(\boldsymbol{A})^{(T-1)/2}\sqrt{T\|\bar{\boldsymbol{\Gamma}}_\infty\|} \leq 2\sqrt{\frac{p}{p+m}}\varepsilon_0/C.$$

This implies $\frac{R_e}{\sigma_u\sqrt{N}} \leq \frac{\varepsilon_0}{2}\sqrt{\frac{Tp+m}{N}}$. Combining this with the $R_z$ bound of Theorem 3.2 and tail bound of (I.6), we obtain

$$\|\boldsymbol{G}^{(\infty)} - \hat{\boldsymbol{G}}^{(\infty)}\| \leq (8\frac{\sigma_z}{\sigma_u} + \varepsilon_0)\sqrt{\frac{Tp+m}{N}},$$

whenever $N$ is stated as above and $T$ obeys $T \geq \max(-\frac{2}{\log(\rho(\boldsymbol{A}))}, T_0)$ where

$$T_0 := 2\frac{c_0 + \log(\Phi(\boldsymbol{A})^2\|\boldsymbol{C}\|\varepsilon_0^{-1}) + \log\alpha}{-\log(\rho(\boldsymbol{A}))} \geq \max(T_1, T_2),$$

and $\alpha$ is a shorthand variable for

$$\alpha = (1 - \rho(\boldsymbol{A})^{1/2})^{-1}\|\boldsymbol{B}\|\sqrt{\frac{N}{Tp+m}} + \sqrt{T\|\bar{\boldsymbol{\Gamma}}_\infty\|(1+m/p)}.$$

Treating $\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}$ related variables in the numerator as constant terms (which are less insightful than the $\log(\rho(\boldsymbol{A}))$ term for our purposes), we find the condition (IV.2).

To proceed, we wish to show the result on Hankel matrices $\boldsymbol{H}^{(\infty)}$ and $\hat{\boldsymbol{H}}^{(\infty)}$. We shall decompose the $\boldsymbol{H}^{(\infty)}$ matrix as $\boldsymbol{H}^{(\infty)} = \boldsymbol{H}_{\text{main}} + \boldsymbol{H}_{\text{tail}}$ (same for $\hat{\boldsymbol{H}}$). $\boldsymbol{H}_{\text{main}}, \hat{\boldsymbol{H}}_{\text{main}}$ are the $m \times p$ blocks corresponding to the first $T$ Markov parameters and their estimates. Observe that $\boldsymbol{H}_{\text{main}}$ lives on the upper-left $T \times T$ submatrix. Furthermore, the set of non-zero blocks in each of its first $T$ block-rows of size $m \times Tp$ is a submatrix of $\boldsymbol{G}$. For instance following Def. 4.1, non-zero rows of $\hat{\boldsymbol{H}}$ are all submatrices of $\hat{\boldsymbol{G}}$. Consequently, using the above bound on $\boldsymbol{G}$ estimate and naively bounding the overall spectral norm of these $T$ nonzero rows (which results in a $\sqrt{T}$ factor), we have that

$$\|\hat{\boldsymbol{H}}_{\text{main}} - \boldsymbol{H}_{\text{main}}\| \leq \sqrt{T} \times \|\boldsymbol{G} - \hat{\boldsymbol{G}}\| \leq (8\frac{\sigma_z}{\sigma_u} + \frac{\varepsilon_0}{2})\sqrt{T \times \frac{Tp+m}{N}},$$

where $\varepsilon_0/2$ instead of $\varepsilon_0$ is due to lack of tail terms. What remains is the $\boldsymbol{H}_{\text{tail}}$ term. Note that $\hat{\boldsymbol{H}}_{\text{tail}} = 0$. $\boldsymbol{H}_{\text{tail}}$ matrix is composed of anti-diagonal blocks that start from $T+1$ till infinity. The non-zero blocks of $i$th anti-diagonal ($i \geq T+1$) are all equal to $\boldsymbol{C}\boldsymbol{A}^{i-2}\boldsymbol{B}$ due to Hankel structure, hence its spectral norm is equal to $\|\boldsymbol{C}\boldsymbol{A}^{i-2}\boldsymbol{B}\|$. Consequently, the spectral norm of $\boldsymbol{H}_{\text{tail}}$ can be obtained by adding the spectral norm of non-zero anti-diagonal matrices which is given by (IV.1) and is upper bounded by $\varepsilon_0/2$ in (I.6). Hence,

$$\|\hat{\boldsymbol{H}}^{(\infty)} - \boldsymbol{H}^{(\infty)}\| \leq \|\hat{\boldsymbol{H}}_{\text{main}} - \boldsymbol{H}_{\text{main}}\| + \|\hat{\boldsymbol{H}}_{\text{tail}} - \boldsymbol{H}_{\text{tail}}\|$$

concluding the proof. ∎

### APPENDIX II
### PROOF OF THE HO-KALMAN ROBUSTNESS

In this section, we provide a proof for the robustness of the Ho-Kalman procedure. Since system is assumed to be observable and controllable and $T_1, T_2$ are assumed to be sufficiently large, rank$(\boldsymbol{L}) = n$ throughout this section. Recall that, given Markov parameter matrices $\boldsymbol{G}, \hat{\boldsymbol{G}}$, the matrices $\boldsymbol{H}, \boldsymbol{H}^-, \boldsymbol{L}, \boldsymbol{H}^+$ (with $\boldsymbol{L} = \boldsymbol{H}^-$ as $\boldsymbol{H}^-$ is rank $n$) correspond to $\boldsymbol{G}$ and the matrices $\hat{\boldsymbol{H}}, \hat{\boldsymbol{H}}^-, \hat{\boldsymbol{L}}, \hat{\boldsymbol{H}}^+$ correspond to $\hat{\boldsymbol{G}}$. We will show that Ho-Kalman state-space realizations corresponding to $\boldsymbol{G}$ and $\hat{\boldsymbol{G}}$ are close to each other as a function of $\|\boldsymbol{G} - \hat{\boldsymbol{G}}\|$. We first provide a proof of Lemma 5.1.

## A. Proof of Lemma 5.1

We wish to show that $H - \hat{H}$ and $L - \hat{L}$ can be upper bounded in terms of $G - \hat{G}$ via (V.1). $H^- - \hat{H}^-$ is a submatrix of $H - \hat{H}$ hence we have

$$\|H^- - \hat{H}^-\| \le \|H - \hat{H}\|.$$

Denote the $i$th block row of $H$ by $H[i]$. Since $H[i]$ (for all $i$) is a submatrix of the Markov parameter matrix $G$, we have that $\|H[i] - \hat{H}[i]\| \le \|G - \hat{G}\|$. Hence, the overall matrix $H$ satisfies

$$\|H - \hat{H}\| = \left\| \begin{bmatrix} H[1] - \hat{H}[1] \\ \vdots \\ H[T_1] - \hat{H}[T_1] \end{bmatrix} \right\| \le \sqrt{T_1} \|G - \hat{G}\|.$$

Similarly, columns of $H$ are also submatrices of $G$. Repeating same argument for columns, yields

$$\|H - \hat{H}\| \le \sqrt{T_2 + 1} \|G - \hat{G}\|.$$

Combining both, we find (V.1). The bound (V.2) is based on singular value perturbation. First, noticing that rows/columns of $H^-$ are again copied from $G$ and carrying out the same argument, we have that

$$\|H^- - \hat{H}^-\| \le \sqrt{\min\{T_1, T_2\}} \|G - \hat{G}\|.$$

Recall that $L = H^-$ and $\hat{L}$ is the rank-$n$ approximations of $\hat{H}^-$. Denoting $i$th singular value of $\hat{H}^-$ by $\sigma_i(\hat{H}^-)$, standard singular value perturbation bound yields

$$\sigma_{n+1}(\hat{H}^-) = \|\hat{H}^- - \hat{L}\| \le \|\hat{H}^- - H^-\|.$$

Consequently, using $L = H^-$,

$$\begin{aligned} \|L - \hat{L}\| &\le \|H^- - \hat{H}^-\| + \|\hat{H}^- - \hat{L}\| \\ &\le 2\|H^- - \hat{H}^-\| \le 2\sqrt{\min\{T_1, T_2\}} \|G - \hat{G}\|. \end{aligned}$$

## B. Robustness of Singular Value Decomposition

The next theorem shows robustness of singular value decompositions of $L$ and $\hat{L}$ in terms of $\|L - \hat{L}\|$. It is obtained by using Lemma 5.14 of [38] and provides simultaneous control over left and right singular vector subspaces. This is essentially similar to results of Wedin and Davis-Kahan [40], [41] with the added advantage of simultaneous control which we crucially need for our result.

*Lemma 2.1:* Suppose $\sigma_n(L) \ge 2\|L - \hat{L}\|$ where $\sigma_n(L)$ is the smallest nonzero singular value (i.e. $n$th largest singular value) of $L$. Let rank $n$ matrices $L$, $\hat{L}$ have singular value decompositions $U\Sigma V^*$ and $\hat{U}\hat{\Sigma}\hat{V}^*$. There exists an $n \times n$ unitary matrix $T$ so that

$$\|U\Sigma^{1/2} - \hat{U}\hat{\Sigma}^{1/2}T\|_F^2 + \|V\Sigma^{1/2} - \hat{V}\hat{\Sigma}^{1/2}T\|_F^2 \le \frac{10n\|L - \hat{L}\|^2}{\sigma_n(L)}.$$

*Proof:* Direct application of Lemma 5.14 of [38] guarantees the existence of a unitary $T$ such that

$$\begin{aligned} \text{LHS} &= \|U\Sigma^{1/2} - \hat{U}\hat{\Sigma}^{1/2}T\|_F^2 + \|V\Sigma^{1/2} - \hat{V}\hat{\Sigma}^{1/2}T\|_F^2 \\ &\le \frac{2}{\sqrt{2}-1} \frac{\|L - \hat{L}\|_F^2}{\sigma_n(L)}. \end{aligned}$$

To proceed, use the fact that $\text{rank}(L - \hat{L}) \le 2n$ to bound $\|L - \hat{L}\|_F \le \sqrt{2n}\|L - \hat{L}\|$. ∎

Observe that our control over the subspace deviation improves as the perturbation $\|L - \hat{L}\|$ gets smaller. The next lemma is a standard result on singular value deviation.

*Lemma 2.2:* Suppose $\sigma_n(L) \ge 2\|L - \hat{L}\|$. Then, $\|\hat{L}\| \le 2\|L\|$ and $\sigma_n(\hat{L}) \ge \sigma_n(L)/2$.

Using these, we will prove the robustness of Ho-Kalman. The robustness will be up to a unitary transformation similar to Lemma 2.1.

## C. Proof of Theorem 5.2

*Proof:* Consider the SVD of $L$ given by $U\Sigma V$ and SVD of $\hat{L}$ given by $\hat{U}\hat{\Sigma}\hat{V}$ where $\Sigma, \hat{\Sigma} \in \mathbb{R}^{n \times n}$ (recall that $\text{rank}(L) = n$ since we assumed system is observable and controllable). Define the observability/controllability matrices $(O = U\Sigma^{1/2}, Q = \Sigma^{1/2}V)$ associated to $H$ and $(\hat{O} = \hat{U}\hat{\Sigma}^{1/2}, \hat{Q} = \hat{\Sigma}^{1/2}\hat{V})$ associated to $\hat{H}$. Lemma 2.1 automatically gives control over these as it states the existence of a unitary matrix $T$ such that

$$\|O - \hat{O}T\|_F^2 + \|Q - T^*\hat{Q}\|_F^2 \le 10n\|L - \hat{L}\|^2/\sigma_n(L).$$

Since $\bar{C}$ is a submatrix of $O$ and $\bar{B}$ is a submatrix of $Q$, we immediately have the same upper bound on $(\bar{C}, \hat{C})$ and $(\bar{B}, \hat{B})$ pairs.

The remaining task is to show that $\hat{A}$ and $\bar{A}$ are close. Let $X = \hat{O}T$, $Y = T^*\hat{Q}$. Now, note that

$$\begin{aligned} \|\bar{A} - T^*\hat{A}T\|_F &= \|O^\dagger H^+ Q^\dagger - T^*\hat{O}^\dagger \hat{H}^+ \hat{Q}^\dagger T\|_F \\ &= \|O^\dagger H^+ Q^\dagger - X^\dagger \hat{H}^+ Y^\dagger\|_F. \quad \text{(II.1)} \end{aligned}$$

Consequently, we can decompose the right hand side as

$$\begin{aligned} \|O^\dagger H^+ Q^\dagger - X^\dagger \hat{H}^+ Y^\dagger\|_F \le &\|(O^\dagger - X^\dagger)H^+ Q^\dagger\|_F + \quad \text{(II.2)} \\ &\|X^\dagger(H^+ - \hat{H}^+)Q^\dagger\|_F + \|X^\dagger \hat{H}^+(Q^\dagger - Y^\dagger)\|_F. \end{aligned}$$

We treat the terms on the right hand side individually. First, pseudo-inverse satisfies the perturbation bound [42], [43]

$$\begin{aligned} \|O^\dagger - X^\dagger\|_F &\le \|O - X\|_F \max\{\|X^\dagger\|^2, \|O^\dagger\|^2\} \\ &\le \sqrt{\frac{10n\|L - \hat{L}\|^2}{\sigma_n(L)}} \max\{\|X^\dagger\|^2, \|O^\dagger\|^2\}. \end{aligned}$$

We need to bound the right hand side. Luckily, Lemma 2.2 trivially yields the control over the top singular values of pseudo-inverses namely

$$\max\{\|X^\dagger\|^2, \|O^\dagger\|^2\} = \max\left\{\frac{1}{\sigma_n(L)}, \frac{1}{\sigma_n(\hat{L})}\right\} \le \frac{2}{\sigma_n(L)}.$$

Combining the last two bounds, we find

$$\|O^\dagger - X^\dagger\|_F \le \frac{2\sqrt{\frac{10n\|L - \hat{L}\|^2}{\sigma_n(L)}}}{\sigma_n(L)}$$

The identical bounds hold for $Q, Y$. For the second term on the right hand side of (II.2), we shall use the estimate

$$\frac{\|X^\dagger(H^+ - \hat{H}^+)Q^\dagger\|_F}{\sqrt{n}} \le \|X^\dagger(H^+ - \hat{H}^+)Q^\dagger\| \le \frac{2\|H^+ - \hat{H}^+\|}{\sigma_n(L)}.$$

Finally, we will use the standard triangle inequality to address the $\hat{H}^+$ term: $\|\hat{H}^+\| \leq \|H^+\| + \|H^+ - \hat{H}^+\|$. Combining all of these, we obtain the following bounds

$$\|(O^\dagger - X^\dagger)H^+ Q^\dagger\|_F \leq \|O^\dagger - X^\dagger\|_F \|H^+\| \|Q^\dagger\|$$
$$\leq \frac{\sqrt{40n\|L - \hat{L}\|^2}}{\sigma_n(L)^{3/2}} \sqrt{\frac{2}{\sigma_n(L)}} \|H^+\|$$
$$\leq \frac{9\sqrt{n}\|L - \hat{L}\|}{\sigma_n^2(L)} \|H^+\|$$

$$\|X^\dagger \hat{H}^+ (Q^\dagger - Y^\dagger)\|_F \leq \|X^\dagger\| \|\hat{H}^+\| \|Q^\dagger - Y^\dagger\|_F$$
$$\leq \frac{9\sqrt{n}\|L - \hat{L}\|}{\sigma_n^2(L)} (\|H^+\| + \|H^+ - \hat{H}^+\|)$$

$$\|X^\dagger (H^+ - \hat{H}^+) Q^\dagger\|_F \leq \frac{2\sqrt{n}\|H^+ - \hat{H}^+\|}{\sigma_n(L)}.$$

Combining these three individual bounds and substituting in (II.2), we find the overall bound (V.5). ∎

### D. Proof of Theorem 5.3

Under provided assumptions, Theorem 3.1 yields $\|G - \hat{G}\| \leq C\sqrt{N_0/N}\log(Nq)$ with high probability for some constant $C > 0$. This yields $\|G - \hat{G}\| \leq \sigma_n(L)/4\sqrt{T}$ under the bound (V.6). This also implies $\|H - \hat{H}\| \leq \sqrt{T}\|G - \hat{G}\| \leq \sigma_n(L)/4$ which in turn implies the condition (V.3). Consequently, inequalities (V.4) and (V.5) of Theorem 5.2 hold. (II.3) follows by using $\|L - \hat{L}\| \leq 2\|H - \hat{H}\|$. This yields

$$\max\{\|\bar{C} - \hat{C}T\|_F, \|O - \hat{O}T\|_F, \|\bar{B} - T^*\hat{B}\|_F, \|Q - T^*\hat{Q}\|_F\}$$
$$\leq 7\sqrt{n}\|H - \hat{H}\|/\sqrt{\sigma_n(L)}. \quad (II.3)$$

Plugging in the bound on $\|H - \hat{H}\|$ leads to (II.3). The result on $\bar{A}$ is slightly more intricate. First, since $H^+$ is a submatrix of $H$

$$\|H^+ - \hat{H}^+\| \leq \|H - \hat{H}\|, \quad \|H^+\| \leq \|H\|$$

Combining this with (V.2), the right hand side of (V.5) can be upper bounded by

$$\text{RHS} = \frac{18\sqrt{n}}{\sigma_n(L)}\left(\frac{\|H - \hat{H}\|}{\sigma_n(L)}\|H\| + \|H - \hat{H}\|\right) \quad (II.4)$$
$$= \frac{18\sqrt{n}\|H - \hat{H}\|}{\sigma_n(L)}\left(\frac{\|H\|}{\sigma_n(L)} + 1\right) \leq \frac{36\sqrt{n}\|H - \hat{H}\|\|H\|}{\sigma_n^2(L)}.$$

Plugging in the finite sample bound on $\|H - \hat{H}\|$ yields (V.8).

## APPENDIX III
## RESTRICTED ISOMETRY OF PARTIAL CIRCULANT MATRICES

To proceed, let us describe the goal of this section. First, we would like to show that $U \in \mathbb{R}^{N \times Tp}$ is well conditioned when $N \gtrsim \mathcal{O}(Tp)$ to ensure least-squares is robust. Next, we would like to have an accurate upper bound on the spectral norm of $U^*W$ to control the impact of noise $w_t$. In particular, we will show that

$$\|U^*W\| \lesssim \sigma_u \sigma_w \sqrt{NT(p+n)}.$$

Both of these goals will be achieved by embedding $U$ and $W$ into proper circulant matrices. The same argument will apply to both scenarios. The key technical tool in our analysis will be the results of Krahmer et al. [36] on restricted isometries of random circulant matrices.

The following theorem is a restatement of Theorem 4.1 of Krahmer et al [36]. We added a minor modification to account for the regime *restricted isometry constant* is greater than 1. This theorem shows that arbitrary submatrices of random circulant matrices are well conditioned. It will play a crucial role in establishing the joint relation of the data matrix $U$ and noise matrix $W$. Main result of [36] characterizes a uniform bound on all submatrices; however we only need a single submatrix for our results. Hence, some of the logarithmic factors below might actually be redundant for the bound we are seeking.

*Theorem 3.1:* Let $C \in \mathbb{R}^{d \times d}$ be a circulant matrix where the first row is distributed as $\mathcal{N}(0, I_d)$. Given $s, d \geq 2$, set $m_0 = c_0 s \log^2(s) \log^2(d)$ for some absolute constant $c_0 > 0$. Pick an $m \times s$ submatrix $S$ of $C$. With probability at least $1 - d^{-\log(d)\log^2(s)}$, $S$ satisfies

$$\left\|\frac{1}{m}S^*S - I\right\| \leq \max\left\{\sqrt{\frac{m_0}{m}}, \frac{m_0}{m}\right\}.$$

This result is proven in the next subsection.

### A. Proof of Theorem 3.1

This proof is a slight modification of the proof of Theorem 4.1 of Krahmer et al. [36] and we will directly borrow their notation and estimates. First, we restate their Theorem 3.1.

*Theorem 3.2:* Let $\mathcal{A}$ be a set of matrices and let $\xi$ be a random vector whose entries $\xi_j$ are standard normal. Let $d_F, d_{2\to2}$ be the Frobenius and spectral norm distance metrics respectively, and understood as the radius when applied to a set. Let $\gamma_2(\cdot)$ be Talagrand's $\gamma$-functional [44] that arises in the study of suprema of random processes and is used to capture the richness of a set. Set

$$E = \gamma_2(\mathcal{A}, \|\cdot\|)(\gamma_2(\mathcal{A}, \|\cdot\|) + d_F(\mathcal{A})) + d_F(\mathcal{A})d_{2\to2}(\mathcal{A}),$$
$$V = d_{2\to2}(\mathcal{A})(\gamma_2(\mathcal{A}, \|\cdot\|) + d_F(\mathcal{A})), \text{ and } U = d_{2\to2}^2(\mathcal{A}).$$

Then, for some absolute constants $c_1, c_2 > 0$ and for all $t > 0$,

$$\mathbb{P}\left(\sup_{A \in \mathcal{A}} |\|A\xi\|_{\ell_2}^2 - \mathbb{E}\|A\xi\|_{\ell_2}^2| \geq c_1 E + t\right) \leq \exp\left(-c_2 \min\left\{\frac{t^2}{V^2}, \frac{t}{U}\right\}\right).$$

Theorem 3.1 is a variation of Theorem 4.1 of [36]. In light of Theorem 3.2, we simply need to adapt the estimates developed during the proof of Theorem 4.1 of [36] for our purposes. We are interested in a fixed submatrix of size $m \times s$ compared to all $s$-column submatrices for fixed $m$-rows. This makes our set $\mathcal{A}$ a subset of their set and also makes their estimates an upper bound on our estimates. Following arguments of [36], for some constant $c_3 > 0$, we have $d_F(\mathcal{A}) = 1$ and

$$d_{2\to2}(\mathcal{A}) \leq \sqrt{s/m}, \quad \gamma_2(\mathcal{A}, \|\cdot\|) \leq c_3\sqrt{s/m}\log(s)\log(d).$$

To proceed, we will apply Theorem 3.2. This will be done in two scenarios depending on whether the isometry constant obeys $\delta \leq 1$ or not. Recall that $m_0 = c_0 s \log^2(s) \log^2(d)$. Below, we pick $c_0$ sufficiently large to compensate for $c_1, c_2, c_3$.

$m \geq m_0$ **case:** We have that $\gamma_2(\mathcal{A}, \| \cdot \|) \leq c_3\sqrt{s/m}\log(s)\log(d) \leq 1$ so that $E \leq \sqrt{s/m} + 2\gamma_2(\mathcal{A}, \|\cdot\|) \leq 3c_3\sqrt{s/m}\log(s)\log(d)$. Similarly, $V \leq 2\sqrt{s/m}$ and $U \leq s/m$. In this case, picking large $c_0$, observe that $c_1 E \leq \sqrt{m_0/(4m)}$. With this, we can pick $t = \sqrt{m_0/(4m)}$ to guarantee $c_1 E + t \leq \sqrt{m_0/m}$. We have that $t^2/V^2 \geq m_0/(16s)$, $t/U \geq t^2/U \geq m_0/(4s)$. Picking $c_0 \geq 16/c_2$, we conclude with the desired probability $\exp(-\log^2(d)\log^2(s))$.

$m < m_0$ **case:** In this case, we have

$$\gamma_2(\mathcal{A}, \| \cdot \|) + d_F(\mathcal{A}) \leq c'\sqrt{s/m}\log(s)\log(d),$$

where $c' = c_3 + \sqrt{c_0}$. Hence, we find

$$E \leq c'c_3(s/m)\log(s)^2\log(d)^2 + \sqrt{s/m}.$$

Observe that, we can ensure i) $c_1\sqrt{s/m} \leq \sqrt{m_0/m}/4 \leq m_0/(4m)$ and ii) $c_1c'c_3(s/m)\log^2(s)\log^2(d) \leq m_0/(4m)$ for sufficiently large constant $c_0$. The latter one follows from the fact that $c'$ grows proportional to $\sqrt{c_0}$ whereas $m_0$ grows proportional to $c_0$. With this, we can pick $t = m_0/(2m)$ which guarantees $c_1 E + t \leq \frac{m_0}{m}$.

To find the probability, we again pick $c_0$ to be sufficiently large to guarantee that i) $c_2 t/U \geq c_2(m_0/(2m))/(s/m) \geq \log^2(s)\log^2(d)$ and ii)

$$\begin{aligned} t^2/V^2 &\geq \frac{(m_0/m)^2}{4(s/m)(c'\sqrt{s/m}\log(s)\log(d))^2} \\ &= \frac{c_0^2 s^2 \log^4(s)\log^4(d)}{4s^2(c')^2\log^2(s)\log^2(d)} \\ &= \frac{c_0^2 \log^2(s)\log^2(d)}{4(c_3 + \sqrt{c_0})^2} \geq \log^2(s)\log^2(d)/c_2, \end{aligned}$$

which concludes the proof by yielding $\exp(-\log^2(d)\log^2(s))$ probability of success.

### B. Conditioning of the Data Matrix

This section and the next one address the minimum singular value of the $U$ matrix and upper bounding the maximum singular value of the $U^*W$ matrix by utilizing Theorem 3.1.

*Lemma 3.3:* Let $U \in \mathbb{R}^{N \times Tp}$ be the input data matrix as described in Section II-A. Suppose the sample size obeys $N \geq cTp\log^2(Tp)\log^2(\bar{N}p)$ for sufficiently large constant $c > 0$. Then, with probability at least $1 - (\bar{N}p)^{-\log^2(Tp)\log(\bar{N}p)}$,

$$2N\sigma_u^2 \geq U^*U \geq N\sigma_u^2/2.$$

*Proof:* The proof will be accomplished by embedding $U$ inside a proper circulant matrix. Let $r(v): \mathbb{R}^d \to \mathbb{R}^d$ be the circulant shift operator which maps a vector $v \in \mathbb{R}^d$ to its single entry circular rotation to the right i.e. $r(v) = [v_d\ v_1\ \ldots\ v_{d-1}] \in \mathbb{R}^d$. Let $C \in \mathbb{R}^{\bar{N}p \times \bar{N}p}$ be a circulant matrix where the first row (transposed) is given by

$$c_1 = [u_{\bar{N}p}^*\ u_{\bar{N}p-1}^*\ \ldots\ u_2^*\ u_1^*]^*.$$

The $i$th row of $C$ is $c_i = r^{i-1}(c_1)$ for $1 \leq i \leq \bar{N}p$. Observe that $C$ is a circulant matrix by construction. For instance all of its diagonal entries are equal to $u_{\bar{N}p,1}$. Additionally, note that second row of $C$ starts with the last entry of $u_1$ hence entries

of $u_i$ do not necessarily lie next to each other. Focusing on the rightmost $Tp$ columns, let $R_{Tp}$ be the operator that returns rightmost $Tp$ entries of a vector. Our first observation is that

$$R_{Tp}(c_1) = \bar{u}_T = [u_T^*\ u_{T-1}^*\ \ldots\ u_2^*\ u_1^*]^*.$$

Secondly, observe that for each $0 \leq i \leq N - 1$

$$R_{Tp}(c_{1+ip}) = [u_{T+i}^*\ u_{T-1+i}^*\ \ldots\ u_{2+i}^*\ u_{1+i}^*]^* = \bar{u}_{T+i}.$$

This implies that $\bar{u}_{T+i}$ is embedded inside right-most $Tp$ columns and $1 + ip$'th row of $C$. Similarly, the input data matrix $U \in \mathbb{R}^{N \times Tp}$ is a submatrix of $C$ with column indices $(\bar{N} - T)p + 1$ to $\bar{N}p$ and row indices $1 + ip$ for $0 \leq i \leq N - 1$. Applying Theorem 3.1, setting $N^* = cTp\log^2(Tp)\log^2(\bar{N}p)$, and adjusting for variance $\sigma_u^2$, with probability at least $1 - (\bar{N}p)^{-\log^2(Tp)\log(\bar{N}p)}$, we have

$$2\sigma_u^2 I \succeq N^{-1}U^*U \succeq \frac{\sigma_u^2}{2}I \implies 2N\sigma_u^2 \succeq U^*U \succeq N\sigma_u^2/2,$$

whenever $N \geq N^*$. ∎

### C. Upper Bounding the Contribution of the Process Noise

*Lemma 3.4:* Recall $U, W$ from (II.6) and (III.3) respectively. Let $q = p + n$ and $N^* = cTq\log^2(Tq)\log^2(\bar{N}q)$ where $c > 0$ is an absolute constant. With probability at least $1 - (\bar{N}q)^{-\log^2(Tq)\log(\bar{N}q)}$,

$$\|U^*W\| \leq \sigma_w\sigma_u\max\{\sqrt{N^*N}, N^*\}.$$

*Proof:* The proof is identical to that of Lemma 3.3. Set $q = p + n$. First, we define $m_t = [\sigma_u^{-1}u_t^*\ \sigma_w^{-1}w_t^*]^* \in \mathbb{R}^q$ and $\bar{m}_i = [m_i^*,\ m_{i-1}^*,\ \ldots\ m_{i-T+1}^*]^* \in \mathbb{R}^{Tq}$. We also define the matrix $M = [\bar{m}_T\ \ldots\ \bar{m}_{T+N-1}]^* \in \mathbb{R}^{N \times Tq}$. Observe that by construction, $\sigma_u^{-1}U, \sigma_w^{-1}W$ are submatrices of $M$. In particular, $(\sigma_u\sigma_w)^{-1}U^*W$ is an off-diagonal submatrix of $M^*M$ of size $Tp \times Tn$. This is due to the facts that a) $\sigma_u^{-1}U$ is a submatrix of $M$ characterized by the column indices

$$\{(i - 1)q + j\ |\ 1 \leq i \leq T,\ 1 \leq j \leq p\},$$

and b) $\sigma_w^{-1}W$ lies at the complementary columns. Observe that the spectral norm of $(\sigma_u\sigma_w)^{-1}U^*W$ is bounded by

$$(\sigma_u\sigma_w)^{-1}\|U^*W\| \leq \|M^*M - NI\|. \qquad (III.1)$$

This is because $(\sigma_u\sigma_w)^{-1}U^*W$ is an off-diagonal submatrix of $M^*M$, it is also a submatrix of $M^*M - kI$ for any $k \in \mathbb{R}$. Spectral norm of a submatrix is upper bounded by the norm of the original matrix hence the claim follows.

In a similar fashion to Lemma 3.3, we complete $M$ to be a full circulant matrix as follows. Let $r(v): \mathbb{R}^d \to \mathbb{R}^d$ be the circulant shift operator as previously. Let $C \in \mathbb{R}^{\bar{N}q \times \bar{N}q}$ be a circulant matrix with first row given by

$$c_1 = [m_{\bar{N}q}^*\ m_{\bar{N}q-1}^*\ \ldots\ m_2^*\ m_1^*]^*.$$

The $i$th row of $C$ is $c_i = r^{i-1}(c_1)$ for $1 \leq i \leq \bar{N}q$. Let $R_{Tq}$ be the operator that returns rightmost $Tq$ entries of a vector. Our first observation is that

$$R_{Tq}(c_1) = \bar{m}_T = [m_T^*\ m_{T-1}^*\ \ldots\ m_2^*\ m_1^*]^*.$$

Secondly, observe that for each $0 \leq i \leq N - 1$

$$\boldsymbol{R}_{Tq}(\boldsymbol{c}_{1+iq}) = [\boldsymbol{m}_{T+i}^* \ \boldsymbol{m}_{T-1+i}^* \ \ldots \ \boldsymbol{m}_{2+i}^* \ \boldsymbol{m}_{1+i}^*]^* = \bar{\boldsymbol{m}}_{T+i}.$$

This implies that $\bar{\boldsymbol{m}}_i$'s are embedded inside the rows of $\boldsymbol{R}_{Tq}(\boldsymbol{C})$ in an equally spaced manner with spacing $q$ for $T \leq i \leq T + N - 1 = \bar{N}$. Hence, $\boldsymbol{M}$ is a $N \times Tq$ submatrix of $\boldsymbol{C}$ where the column indices are the last $Tq$ columns and the row indices are $1, 1 + q, \ldots, 1 + (N - 1)q$.

With this observation, by Theorem 3.1, we have, for

$$N^* = cTq \log^2(Tq) \log^2(\bar{N}q),$$

with probability at least $1 - (\bar{N}q)^{-\log^2(Tq)\log(\bar{N}q)}$,

$$\left\| \frac{1}{N} \boldsymbol{M}^* \boldsymbol{M} - \boldsymbol{I} \right\| \leq \max\left\{ \sqrt{\frac{N^*}{N}}, \frac{N^*}{N} \right\},$$

which in turn implies $\|\boldsymbol{U}^* \boldsymbol{W}\| \leq \sigma_w \sigma_u \max\{\sqrt{N^* N}, N^*\}$ via inequality (III.1). ∎

## APPENDIX IV
## BOUNDING THE ERROR DUE TO THE UNKNOWN STATE

The goal of this section is bounding the estimation error due to the $\boldsymbol{e}_t = \boldsymbol{C} \boldsymbol{A}^{T-1} \boldsymbol{x}_{t-T+1}$ term. As described in Section II-A and (III.3), we form the matrices $\boldsymbol{E} = [\boldsymbol{e}_T \ \ldots \ \boldsymbol{e}_{\bar{N}}]^*$ and $\boldsymbol{U} = [\bar{\boldsymbol{u}}_T \ \ldots \ \bar{\boldsymbol{u}}_{\bar{N}}]^*$. Our interest in this section is bounding $\|\boldsymbol{U}^* \boldsymbol{E}\|$. This term captures the impact of approximating the system with a finite impulse response of length $T$. We will show that

$$\|\boldsymbol{U}^* \boldsymbol{E}\| \lesssim \sigma_u \sqrt{(Tp + m)NT \|\boldsymbol{\Gamma}_\infty\| \|\boldsymbol{C}\boldsymbol{A}^{T-1}\|^2}.$$

The main challenge in analyzing $\boldsymbol{U}^* \boldsymbol{E}$ is the fact that $\{\boldsymbol{e}_t\}_{t=T}^{\bar{N}}$ terms and $\{\bar{\boldsymbol{u}}_t\}_{t=T}^{\bar{N}}$ terms are dependent. In fact $\boldsymbol{e}_t$ contains a $\boldsymbol{u}_\tau$ component inside for any $\tau \leq t - T$. The following theorem is our main result on bounding this term which carefully addresses these dependencies.

*Theorem 4.1:* Suppose we are given $\boldsymbol{U}, \boldsymbol{E}$, as described in Section II-A and (III.3). Define

$$\gamma = \frac{\|\boldsymbol{\Gamma}_\infty\| \Phi(\boldsymbol{A})^2 \|\boldsymbol{C}\boldsymbol{A}^{T-1}\|^2}{1 - \rho(\boldsymbol{A})^T} \tag{IV.1}$$

and suppose $N \geq T$. Then, with probability at least $1 - T(\exp(-100Tp) + 2\exp(-100m))$,

$$\|\boldsymbol{U}^* \boldsymbol{E}\| \leq c\sigma_u \sqrt{T \max\{N, \frac{mT}{1 - \rho(\boldsymbol{A})^{T/2}}\} \max\{Tp, m\}\gamma}.$$

*Proof:* We first decompose $\boldsymbol{U}^* \boldsymbol{E} = \sum_{t=T}^{\bar{N}} \bar{\boldsymbol{u}}_t \boldsymbol{e}_t^*$ into sum of $T$ smaller products. Given $0 \leq t < T$, create sequences $S_t = \{t+T, t+2T, \ldots, t+N_tT\}$ where $N_t$ is the largest integer satisfying $t + N_t T \leq \bar{N}$. Each sequence has length $N_t$ which is at least $\lfloor N/T \rfloor$ and at most $\lfloor N/T \rfloor + 1$. With this, we form the matrices

$$\boldsymbol{U}_t = [\bar{\boldsymbol{u}}_{t+T}, \ \bar{\boldsymbol{u}}_{t+2T}, \ \ldots, \ \bar{\boldsymbol{u}}_{t+N_tT}]^*,$$
$$\boldsymbol{E}_t = [\boldsymbol{e}_{t+T}, \ \boldsymbol{e}_{t+2T}, \ \ldots, \ \boldsymbol{e}_{t+N_tT}]^*. \tag{IV.2}$$

Then, $\boldsymbol{U}^* \boldsymbol{E}$ can be decomposed as

$$\boldsymbol{U}^* \boldsymbol{E} = \sum_{t=0}^{T-1} \boldsymbol{U}_t^* \boldsymbol{E}_t \implies \|\boldsymbol{U}^* \boldsymbol{E}\| \leq \sum_{t=0}^{T-1} \|\boldsymbol{U}_t^* \boldsymbol{E}_t\|. \tag{IV.3}$$

Corollary 4.3 provides a probabilistic spectral norm bound on each term of this decomposition on the right hand side. In particular, applying Corollary 4.3, substituting $\upsilon$ definition, and union bounding over $T$ terms, for all $t$, we obtain

$$\|\boldsymbol{U}_t^* \boldsymbol{E}_t\| \leq c\sigma_u \sqrt{\max\{N, \frac{mT}{1 - \rho(\boldsymbol{A})^{T/2}}\} \max\{p, m/T\}\gamma},$$

with probability at least $1 - T(\exp(-Tq) + 2\exp(-100m))$. This gives the advertised bound on $\boldsymbol{U}^* \boldsymbol{E}$ via (IV.3). ∎

### A. Upper Bounding the Components of the Unknown State Decomposition

Our goal in this section is providing an upper bound on the spectral norm of $\boldsymbol{U}_t^* \boldsymbol{E}_t$ which is described in (IV.2). The following lemma provides a bound that decays with $1/\sqrt{N_t}$. The main tools in our analysis are the probabilistic upper bound on the $\boldsymbol{E}_t$ matrix developed in Section IV-B and martingale concentration bound that was developed and utilized by the recent work of Simchowitz et al [9]. Below we state our bound in the more practical setup $m \leq n$ to avoid redundant notation. In general, our bound scales with $\min\{m, n\}$.

*Theorem 4.2:* Define $\gamma = \frac{\|\boldsymbol{\Gamma}_\infty\| \Phi(\boldsymbol{A})^2 \|\boldsymbol{C}\boldsymbol{A}^{T-1}\|^2}{1 - \rho(\boldsymbol{A})^T}$. $\boldsymbol{U}_t^* \boldsymbol{E}_t$ obeys

$$\|\boldsymbol{U}_t^* \boldsymbol{E}_t\| \leq c_0 \sigma_u \sqrt{\tau \max\{Tp, m\} N_t \gamma},$$

with probability at least $1 - \exp(-100 \max\{Tp, m\}) - 2\exp(-c\tau N_t(1 - \rho(\boldsymbol{A})^{T/2}) + 3m)$ for $\tau \geq 1$.

*Proof:* Given matrices $\boldsymbol{U}_t, \boldsymbol{E}_t$, define the filtrations $\mathcal{F}_i = \sigma(\{\boldsymbol{u}_j, \boldsymbol{w}_j\}_{j=1}^{t+iT})$ for $1 \leq i \leq N_t$. According to this definition $\bar{\boldsymbol{u}}_{t+iT}$ is independent of $\mathcal{F}_{i-1}$ and $\bar{\boldsymbol{u}}_{t+iT} \in \mathcal{F}_i$. The reason is earliest input vector contained by $\bar{\boldsymbol{u}}_{t+iT}$ has index $t+1+(i-1)T$ which is larger than $t + (i-1)T$. Additionally, observe that $\boldsymbol{e}_{t+iT} \in \mathcal{F}_{i-1}$ as $\boldsymbol{e}_{t+iT}$ is a deterministic function of $\boldsymbol{x}_{t+1+(i-1)T}$ which is a function of $\{\boldsymbol{u}_j, \boldsymbol{w}_j\}_{j=1}^{t+(i-1)T}$.

We would like to use the fact that, for each $i$, $\boldsymbol{e}_{t+iT}$ and $\bar{\boldsymbol{u}}_{t+iT}$ are independent. Let $\boldsymbol{X}_t = [\boldsymbol{x}_{t+1} \ \ldots \ \boldsymbol{x}_{t+1+(N_t-1)T}]^*$ so that $\boldsymbol{E}_t = \boldsymbol{X}_t(\boldsymbol{C}\boldsymbol{A}^{T-1})^*$. In light of Lemma 4.5, we will use a covering bound on the matrix

$$\boldsymbol{U}_t^* \boldsymbol{E}_t = \boldsymbol{U}_t^* \boldsymbol{X}_t(\boldsymbol{C}\boldsymbol{A}^{T-1})^*.$$

Let $\mathcal{C}_1$ be a $1/4$ $\ell_2$-cover of the unit sphere $\mathcal{S}^{Tp-1}$ and $\mathcal{C}_2$ be a $1/4$ $\ell_2$-cover of the unit sphere in the row space of $\boldsymbol{C}$. There exists such covers satisfying $\log|\mathcal{C}_1| \leq 3Tp$ and $\log|\mathcal{C}_2| \leq 3\min\{m, n\} \leq 3m$. Pick $\boldsymbol{a}, \boldsymbol{b}$ from $\mathcal{C}_1, \mathcal{C}_2$ respectively. Let $W_i = \boldsymbol{a}^* \bar{\boldsymbol{u}}_{t+iT}$ and $Z_i = \boldsymbol{b}^* \boldsymbol{e}_{t+iT}$. Observe that

$$\sum_{i=1}^{N_t} W_i Z_i = \boldsymbol{a}^* (\boldsymbol{U}_t^* \boldsymbol{E}_t) \boldsymbol{b}.$$

We next show that $\sum_{i=1}^{N_t} W_i Z_i$ is small with high probability. Applying Lemma 4.6, we find that, for $\tau \geq 2$, with probability at least $1 - 2\exp(-c\tau N_t(1 - \rho(\boldsymbol{A})^{T/2}))$,

$$\|\boldsymbol{E}_t \boldsymbol{b}\|_{\ell_2}^2 = \sum_{i=1}^{N_t} Z_i^2 \leq \tau N_t \gamma, \tag{IV.4}$$

where our definition of $\gamma$ accounts for the $\|\boldsymbol{\Gamma}_\infty\|$ factor. We will use this bound to ensure Lemma 4.4 is applicable with high probability. Since $\bar{\boldsymbol{u}}_{t+iT}$ has $\mathcal{N}(0, \sigma_u^2)$ entries, applying Lemma 4.4, we obtain

This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication. Citation information: DOI 10.1109/TAC.2021.3083651, IEEE Transactions on Automatic Control

AUTHOR *et al.*: PREPARATION OF PAPERS FOR IEEE TRANSACTIONS AND JOURNALS (FEBRUARY 2017) 15

$$\mathbb{P}(\{\sum_{i=1}^{N_t} W_i Z_i \geq t\} \bigcap \{\sum_{i=1}^{N_t} Z_i^2 \leq \tau N_t \gamma\}) \leq \exp(-\frac{t^2}{c\tau\sigma_u^2 N_t \gamma}).$$

for some absolute constant $c > 0$. Picking $t = 11\sigma_u\sqrt{c\tau\max\{Tp, m\}N_t\gamma}$, we find

$$\mathbb{P}(\{\sum_{i=1}^{N_t} W_i Z_i \geq t\} \bigcap \{\sum_{i=1}^{N_t} Z_i^2 \leq cN_t\gamma\}) \leq e^{-120\max\{Tp, m\}}.$$

Defining variables $W_i(\boldsymbol{a})$ for each $\boldsymbol{a} \in \mathcal{C}_1$, and events $E(\boldsymbol{a}) = \{\sum_{i=1}^{N_t} W_i(\boldsymbol{a})Z_i \geq t\}$, applying a union bound, we obtain,

$$\mathbb{P}(\{\bigcup_{\boldsymbol{a} \in \mathcal{C}_1} E(\boldsymbol{a})\} \bigcap \{\sum_{i=1}^{N_t} Z_i^2 \leq cN_t\gamma\}) \leq \exp(-110\max\{Tp, m\}).$$

Combining this bound with (IV.4), we find that, for a *fixed $\boldsymbol{b}$ and for all $\boldsymbol{a}$*, with probability at least $1 - \exp(-110\max\{Tp, m\}) - 2\exp(-c\tau N_t(1 - \rho(\boldsymbol{A})^{T/2}))$, we have

$$\boldsymbol{a}^* \boldsymbol{U}_t^* \boldsymbol{E}_t \boldsymbol{b} = \sum_{i=1}^{N_t} W_i Z_i \leq c_0 \sigma_u \sqrt{\tau \max\{Tp, m\}N_t\gamma}, \quad \text{(IV.5)}$$

for some $c_0 > 0$. Applying a union bound over all $\boldsymbol{b} \in \mathcal{C}_2$, with probability at least $1 - \exp(-100\max\{Tp, m\}) - 2\exp(-c\tau N_t(1 - \rho(\boldsymbol{A})^{T/2}) + 3m)$, we find that (IV.5) holds for all $\boldsymbol{a}, \boldsymbol{b}$. Overall, we found that for all $\boldsymbol{a}, \boldsymbol{b}$ pairs in the 1/4 covers, $\boldsymbol{a}^*(\boldsymbol{U}_t^* \boldsymbol{E}_t)\boldsymbol{b} \leq \kappa = c_0\sigma_u\sqrt{\tau\max\{Tp, m\}N_t\gamma}$. Applying Lemma 4.5, this implies $\|\boldsymbol{U}_t^* \boldsymbol{E}_t\| \leq 2\kappa$. ∎
The following corollary simplifies the result when $N \geq T$ which is the interesting regime for our purposes.

*Corollary 4.3:* Assume $N \geq T$. With probability at least $1 - \exp(-100Tp) - 2\exp(-100m)$, we have $\|\boldsymbol{U}_t^* \boldsymbol{E}_t\| \leq c'\sigma_u\sqrt{\max\{N, \frac{mT}{1-\rho(\boldsymbol{A})^{T/2}}\}\max\{p, m/T\}\gamma}$ for some constant $c' > 0$.

*Proof:* $N \geq T$ implies $N_t \geq \lfloor N/T \rfloor \geq N/(2T)$. In Theorem 4.2, pick $\tau = \max\{1, c_1\frac{mT}{N(1-\rho(\boldsymbol{A})^{T/2})}\}$ for $c_1 = 206/c$. The choice of $\tau$ guarantees the probability exponent $c\tau N_t(1 - \rho(\boldsymbol{A})^{T/2}) - 3m \geq 100m$. To conclude, observe that $c_0\sigma_u\sqrt{\tau\max\{Tp, m\}N_t\gamma} \leq c'\sigma_u\sqrt{\max\{1, \frac{mT}{N(1-\rho(\boldsymbol{A})^{T/2})}\}\max\{p, m/T\}N\gamma}$ for an absolute constant $c' > 0$. ∎
For completeness, we restate the subgaussian Martingale concentration lemma of Simchowitz et al. which is Lemma 4.2 of [9].

*Lemma 4.4:* Let $\{\mathcal{F}_t\}_{t\geq 1}$ be a filtration, $\{Z_t, W_t\}_{t\geq 1}$ be real valued processes adapted to $\mathcal{F}_t, \mathcal{F}_{t+1}$ respectively (i.e. $Z_t \in \mathcal{F}_t, W_t \in \mathcal{F}_{t+1}$). Suppose $W_t \mid \mathcal{F}_t$ is a $\sigma^2$-sub-gaussian random variable with mean zero. Then

$$\mathbb{P}(\{\sum_{t=1}^{T} Z_t W_t \geq \alpha\} \bigcap \{\sum_{t=1}^{T} Z_t^2 \leq \beta\}) \leq \exp(-\frac{\alpha^2}{2\sigma^2\beta})$$

This lemma implies that $\sum_{t=1}^{T} Z_t W_t$ can essentially be treated as an inner product between a deterministic sequence $Z_t$ and an i.i.d. subgaussian sequence $W_t$.

The following lemma is a slight modification of the standard covering arguments.

*Lemma 4.5 (Covering bound):* Given matrices $\boldsymbol{A} \in \mathbb{R}^{n_1 \times N}, \boldsymbol{B} \in \mathbb{R}^{N \times n_2}$, let $\boldsymbol{M} = \boldsymbol{A}\boldsymbol{B}$. Let $\mathcal{C}_1$ be a 1/4-cover of the unit sphere $\mathcal{S}^{n_1-1}$ and $\mathcal{C}_2$ be a 1/4-cover of the unit sphere in the row space of $\boldsymbol{B}$ (which is at most $\min\{N, n_2\}$ dimensional). Suppose for all $\boldsymbol{a} \in \mathcal{C}_1, \boldsymbol{b} \in \mathcal{C}_2$, we have that $\boldsymbol{a}^* \boldsymbol{M} \boldsymbol{b} \leq \gamma$. Then, $\|\boldsymbol{M}\| \leq 2\gamma$.

*Proof:* Pick unit length vectors $\boldsymbol{x}, \boldsymbol{y}$ achieving $\boldsymbol{x}^* \boldsymbol{M} \boldsymbol{y} = \|\boldsymbol{M}\|$. Let $S$ be the row space of $\boldsymbol{B}$. Observe that $\boldsymbol{y} \in S$. Otherwise, its normalized projection on $S$, $\mathcal{P}_S(\boldsymbol{y})/\|\mathcal{P}_S(\boldsymbol{y})\|_{\ell_2}$ achieves a strictly better inner product with $\boldsymbol{x}^* \boldsymbol{M}$. Pick 1/4 close neighbors $\boldsymbol{a}, \boldsymbol{b}$ of $\boldsymbol{x}, \boldsymbol{y}$ from the covers $\mathcal{C}_1, \mathcal{C}_2$. Then,

$$\boldsymbol{x}^* \boldsymbol{M} \boldsymbol{y} = \boldsymbol{a}^* \boldsymbol{M} \boldsymbol{b} + (\boldsymbol{x}-\boldsymbol{a})^* \boldsymbol{M} \boldsymbol{b} + \boldsymbol{x}^* \boldsymbol{M}(\boldsymbol{y}-\boldsymbol{b}) \leq \gamma + \boldsymbol{x}^* \boldsymbol{M} \boldsymbol{y}/2,$$

due to the maximality of $\boldsymbol{x}, \boldsymbol{y}$. This yields $\boldsymbol{x}^* \boldsymbol{M} \boldsymbol{y} \leq 2\gamma$. ∎

### B. Bounding the Inner Products with the Unknown State

In this section, we develop probabilistic upper bounds for the random variable $\boldsymbol{E}_t \boldsymbol{a}$ where $\boldsymbol{a}$ is a fixed vector and $\boldsymbol{E}_t$ is as defined in (IV.2).

*Lemma 4.6:* Let $\boldsymbol{E}_t \in \mathbb{R}^{N_t \times m}$ be the matrix composed of the rows $\boldsymbol{e}_{t+iT} = \boldsymbol{C}\boldsymbol{A}^{T-1}\boldsymbol{x}_{t+1+iT}$. Define

$$\gamma = \frac{\Phi(\boldsymbol{A})^2 \|\boldsymbol{C}\boldsymbol{A}^{T-1}\|^2}{1 - \rho(\boldsymbol{A})^T}.$$

Given a unit length vector $\boldsymbol{a} \in \mathbb{R}^m$, for all $\tau \geq 2$ and for some absolute constant $c > 0$, we have that

$$\mathbb{P}(\|\boldsymbol{E}_t \boldsymbol{a}\|_{\ell_2}^2 \geq \tau N_t \|\boldsymbol{\Gamma}_\infty\|\gamma) \leq 2\exp(-c\tau N_t(1 - \rho(\boldsymbol{A})^{T/2})).$$

*Proof:* Let $\boldsymbol{d}_t = \boldsymbol{x}_t - \boldsymbol{A}^T \boldsymbol{x}_{t-T}$. By construction (i.e. due to the state-space recursion (II.2)), $\boldsymbol{d}_t$ is independent of $\boldsymbol{x}_{t-T}$. We can write $\boldsymbol{x}_{t+iT}$ as

$$\boldsymbol{x}_{t+iT} = \sum_{j=1}^{i} \boldsymbol{A}^{(i-j)T} \boldsymbol{d}_{t+jT} + \boldsymbol{A}^{iT} \boldsymbol{x}_t. \quad \text{(IV.6)}$$

We wish to understand the properties of the random variable $\|\boldsymbol{E}_t \boldsymbol{a}\|_{\ell_2}^2$ which is same as,

$$s_{\boldsymbol{a}} = \sum_{i=1}^{N_t} (\boldsymbol{a}^* \boldsymbol{e}_{t+iT})^2 = \sum_{i=0}^{N_t-1} ((\boldsymbol{a}^* \boldsymbol{C}\boldsymbol{A}^{T-1})\boldsymbol{x}_{t+1+iT})^2.$$

Denote $\bar{\boldsymbol{a}} = (\boldsymbol{C}\boldsymbol{A}^{T-1})^* \boldsymbol{a}$, $\boldsymbol{a}_j = (\boldsymbol{A}^{jT})^* \bar{\boldsymbol{a}}$, $\boldsymbol{g}_0 = \boldsymbol{x}_{t+1}$, and $\boldsymbol{g}_i = \boldsymbol{d}_{t+1+iT}$ for $N_t - 1 \geq i \geq 1$, all of which are $n$ dimensional vectors. Using these change of variables and applying expansion (IV.6), the $i$th component of the sum $s_{\boldsymbol{a}}$ is given by

$$s_{\boldsymbol{a},i} = (\bar{\boldsymbol{a}}^* \boldsymbol{x}_{t+1+iT})^2 = (\bar{\boldsymbol{a}}^* \sum_{j=0}^{i} \boldsymbol{A}^{(i-j)T} \boldsymbol{g}_j)^2$$

$$= (\sum_{j=0}^{i} \boldsymbol{a}_{i-j}^* \boldsymbol{g}_j)^2 = \sum_{0 \leq j,k \leq i} \boldsymbol{a}_{i-j}^* \boldsymbol{g}_j \boldsymbol{a}_{i-k}^* \boldsymbol{g}_k. \quad \text{(IV.7)}$$

Observe that, summing over all $s_{\boldsymbol{a},i}$ for $0 \leq i \leq N_t - 1$, the multiplicative coefficient of the $\boldsymbol{g}_j \boldsymbol{g}_k^*$ pair is given by the matrix,

$$\boldsymbol{M}_{j,k} = \begin{cases} \sum_{N_t > i \geq \max\{j,k\}} \boldsymbol{a}_{i-j}\boldsymbol{a}_{i-k}^* & \text{if } j \neq k, \\ \sum_{N_t > i \geq j} \boldsymbol{a}_{i-j}\boldsymbol{a}_{i-j}^* = \sum_{i=0}^{N_t-1-j} \boldsymbol{a}_i \boldsymbol{a}_i^* & \text{if } j = k \end{cases} \quad \text{(IV.8)}$$

Next, we show that these $\boldsymbol{M}_{j,k}$ submatrices have bounded spectral, Frobenius and nuclear norms (nuclear norm is the sum of the singular values of a matrix). This follows by writing each submatrix as a sum of rank 1 matrices and using the fact

that spectral radius of $\boldsymbol{A}$ is strictly bounded from above by 1.

$$
\begin{aligned}
\|\boldsymbol{M}_{j,k}\| \leq \|\boldsymbol{M}_{j,k}\|_F \leq \|\boldsymbol{M}_{j,k}\|_\star &\leq \sum_{i \geq \max\{j,k\}} \|\boldsymbol{a}_{i-j}\boldsymbol{a}_{i-k}^*\|_\star \\
&= \sum_{i \geq \max\{j,k\}} \|\boldsymbol{a}_{i-j}\boldsymbol{a}_{i-k}^*\| \\
&\leq \sum_{i \geq \max\{j,k\}} \|(\boldsymbol{A}^{(i-j)T})^*\bar{\boldsymbol{a}}\bar{\boldsymbol{a}}^*\boldsymbol{A}^{(i-k)T}\| \\
&\leq \sum_{i \geq \max\{j,k\}} \|\bar{\boldsymbol{a}}\|_{\ell_2}^2 \|\boldsymbol{A}^{(i-j)T}\|\|\boldsymbol{A}^{(i-k)T}\| \\
&\leq \sum_{i=0}^{\infty} \|\bar{\boldsymbol{a}}\|_{\ell_2}^2 \rho(\boldsymbol{A})^{|j-k|T/2}\rho(\boldsymbol{A})^{iT}\Phi(\boldsymbol{A})^2 \\
&\leq \frac{\Phi(\boldsymbol{A})^2\|\bar{\boldsymbol{a}}\|_{\ell_2}^2}{1-\rho(\boldsymbol{A})^T}\rho(\boldsymbol{A})^{|j-k|T/2}.
\end{aligned}
$$

To further simplify, observe that $\|\bar{\boldsymbol{a}}\|_{\ell_2}^2 \leq \|\boldsymbol{C}\boldsymbol{A}^{T-1}\|^2$ as $\|\boldsymbol{a}\|_{\ell_2} = 1$. Setting $\gamma = \frac{\Phi(\boldsymbol{A})^2\|\boldsymbol{C}\boldsymbol{A}^{T-1}\|^2}{1-\rho(\boldsymbol{A})^T}$, we have

$$\|\boldsymbol{M}_{j,k}\|, \|\boldsymbol{M}_{j,k}\|_F, \|\boldsymbol{M}_{j,k}\|_\star \leq \gamma\rho(\boldsymbol{A})^{|j-k|T/2}. \tag{IV.9}$$

Based on the submatrices $\boldsymbol{M}_{j,k}$, create the $N_t n \times N_t n$ matrix $\boldsymbol{M}$. Now we define the vector $\bar{\boldsymbol{g}} = [\boldsymbol{g}_0^* \ \boldsymbol{g}_1^* \ \cdots \ \boldsymbol{g}_{N_t-1}^*]^*$. Observe that, following (IV.7) and (IV.8), by construction,

$$s_{\boldsymbol{a}} = \bar{\boldsymbol{g}}^*\boldsymbol{M}\bar{\boldsymbol{g}} = \sum_{0 \leq j,k < N_t} \boldsymbol{g}_j^*\boldsymbol{M}_{j,k}\boldsymbol{g}_k. \tag{IV.10}$$

This puts $s_{\boldsymbol{a}}$ in a form for which Hanson-Wright Theorem is applicable [45], [46]. To apply Hanson-Wright Theorem, let us first bound the expectation of $s_{\boldsymbol{a}}$. Since $\{\boldsymbol{g}_i\}_{i=0}^{N_t-1}$'s are truncations of the state vector, we have that $\boldsymbol{\Sigma}(\boldsymbol{g}_i) \preceq \boldsymbol{\Sigma}(\boldsymbol{x}_{t+1+iT}) \preceq \boldsymbol{\Gamma}_\infty$. Write $\boldsymbol{g}_i = \boldsymbol{\Sigma}(\boldsymbol{g}_i)^{1/2}\boldsymbol{h}_i$ for some $\boldsymbol{h}_i \sim \mathcal{N}(0, \boldsymbol{I}_n)$. Using independence of $\boldsymbol{h}_i, \boldsymbol{h}_j$ for $i \neq j$ and $\boldsymbol{\Sigma}(\boldsymbol{g}_i) \preceq \boldsymbol{\Gamma}_\infty$, we have that

$$
\begin{aligned}
\mathbb{E}[s_{\boldsymbol{a}}] &= \sum_{i=0}^{N_t-1} \mathbb{E}[\boldsymbol{g}_i^*\boldsymbol{M}_{i,i}\boldsymbol{g}_i] = \sum_{i=0}^{N_t-1} \mathbb{E}[\boldsymbol{h}_i^*\boldsymbol{\Sigma}(\boldsymbol{g}_i)^{1/2}\boldsymbol{M}_{i,i}\boldsymbol{\Sigma}(\boldsymbol{g}_i)^{1/2}\boldsymbol{h}_i] \\
&= \sum_{i=0}^{N_t-1} \operatorname{tr}(\boldsymbol{\Sigma}(\boldsymbol{g}_i)^{1/2}\boldsymbol{M}_{i,i}\boldsymbol{\Sigma}(\boldsymbol{g}_i)^{1/2}) \tag{IV.11} \\
&\leq \sum_{i=0}^{N_t-1} \|\boldsymbol{\Sigma}(\boldsymbol{g}_i)\|\operatorname{tr}(\boldsymbol{M}_{i,i}) \leq \sum_{i=0}^{N_t-1} \|\boldsymbol{\Gamma}_\infty\|\operatorname{tr}(\boldsymbol{M}_{i,i}) \\
&\leq N_t\|\boldsymbol{\Gamma}_\infty\|\gamma. \tag{IV.12}
\end{aligned}
$$

In (IV.11), we utilized the fact that for positive semidefinite matrices trace is equal to the nuclear norm and then we used the fact that nuclear norm of the product obeys $\|\boldsymbol{X}\boldsymbol{Y}\|_\star \leq \|\boldsymbol{X}\|_\star\|\boldsymbol{Y}\|$ [47]. Finally, we upper bounded $\|\boldsymbol{\Sigma}(\boldsymbol{g}_i)\|$ by using the relation $\boldsymbol{\Sigma}(\boldsymbol{g}_i) \preceq \boldsymbol{\Gamma}_\infty$. Bounded $\|\boldsymbol{\Sigma}(\boldsymbol{g}_i)\|$ also implies that the Gaussian vector $\boldsymbol{g}_i$ obeys the "concentration property" (Definition 2.1 of [45]) with $K = \mathcal{O}(\sqrt{\|\boldsymbol{\Gamma}_\infty\|})$ as Lipschitz functions of Gaussians concentrate. Recalling (IV.10), the Hanson-Wright Theorem of [45] states that

$$\mathbb{P}(s_{\boldsymbol{a}} \geq \mathbb{E}[s_{\boldsymbol{a}}]+t) \leq 2\exp(-c\min\{\frac{t^2}{\|\boldsymbol{\Gamma}_\infty\|^2\|\boldsymbol{M}\|_F^2}, \frac{t}{\|\boldsymbol{\Gamma}_\infty\|\|\boldsymbol{M}\|}\}).$$

To proceed, we upper bound $\|\boldsymbol{M}\|_F$ and $\|\boldsymbol{M}\|$. First, recall again that $\|\boldsymbol{M}_{i,j}\|_F \leq \gamma\rho(\boldsymbol{A})^{|i-j|T/2}$. Adding these over all $i, j$

pairs, using (IV.9) and the fact that there are at most $2N_t$ pairs with fixed difference $|i-j| = \tau$, we obtain

$$
\begin{aligned}
\|\boldsymbol{M}\|_F^2 = \sum_{i,j} \|\boldsymbol{M}_{i,j}\|_F^2 &\leq \sum_{0 \leq i,j \leq N_t-1} \gamma^2\rho(\boldsymbol{A})^{|i-j|T} \\
&\leq 2N_t\gamma^2 \sum_{\tau=0}^{N_t-1} \rho(\boldsymbol{A})^{\tau T} \leq \frac{2\gamma^2 N_t}{1-\rho(\boldsymbol{A})^T}.
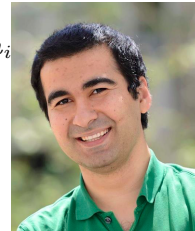\end{aligned}
$$

To assess the spectral norm, we decompose $\boldsymbol{M}$ into $2N_t - 1$ block permutation matrices $\{\boldsymbol{M}^{(i)}\}_{i=-N_t+1}^{N_t-1}$. $\boldsymbol{M}^{(0)}$ is the main diagonal of $\boldsymbol{M}$, and $\boldsymbol{M}^{(i)}$ is the $i$th off-diagonal that contains only the submatrices $\boldsymbol{M}_{j,k}$ with fixed difference $j - k = i$. By construction $\|\boldsymbol{M}^{(i)}\| \leq \gamma\rho(\boldsymbol{A})^{|i|T/2}$ as each nonzero submatrix satisfies the same spectral norm bound. Hence using (IV.9),

$$\|\boldsymbol{M}\| \leq \sum_{i=-N_t+1}^{N_t-1} \|\boldsymbol{M}^{(i)}\| \leq \gamma\left(\frac{2}{1-\rho(\boldsymbol{A})^{T/2}}-1\right) \leq \frac{2\gamma}{1-\rho(\boldsymbol{A})^{T/2}}.$$

With these, setting $t = \tau N_t\|\boldsymbol{\Gamma}_\infty\|\gamma$ and using (IV.12) and bounds on $\|\boldsymbol{M}\|_F, \|\boldsymbol{M}\|$, for $\tau \geq 1$ and using $K = \mathcal{O}(\sqrt{\|\boldsymbol{\Gamma}_\infty\|})$, and applying Theorem 2.3 of [45], we find the concentration bound

$$
\begin{aligned}
\mathbb{P}(s_{\boldsymbol{a}} &\geq (\tau+1)N_t\|\boldsymbol{\Gamma}_\infty\|\gamma) \\
&\leq 2\exp(-2c\tau\min\{\frac{(N_t\|\boldsymbol{\Gamma}_\infty\|\gamma)^2}{\|\boldsymbol{\Gamma}_\infty\|^2\frac{2\gamma^2 N_t}{1-\rho(\boldsymbol{A})^T}}, \frac{N_t\|\boldsymbol{\Gamma}_\infty\|\gamma}{\|\boldsymbol{\Gamma}_\infty\|\frac{2\gamma}{1-\rho(\boldsymbol{A})^{T/2}}}\}) \\
&\leq 2\exp(-c\tau\min\{N_t(1-\rho(\boldsymbol{A})^T), N_t(1-\rho(\boldsymbol{A})^{T/2})\}) \\
&= 2\exp(-c\tau N_t(1-\rho(\boldsymbol{A})^{T/2})),
\end{aligned}
$$

which is the desired result after $1 + \tau \leftrightarrow \tau$ substitution and using the initial assumption of $\tau \geq 2$. ∎

**Samet Oymak** is an assistant professor in the Department of Electrical and Computer Engineering, at the University of California, Riverside. He received his MS and PhD degrees from California Institute of Technology; where he was awarded the Wilts Prize for the best thesis in Electrical Engineering. Before joining UCR, he spent time at Google and financial industry, and prior to that he was a fellow at the Simons Institute and a postdoctoral scholar at UC Berkeley. His research explores the mathematical foundations of data science and machine learning by using tools from optimization and statistics. His research interests include mathematical optimization, reinforcement learning, deep learning theory, and high-dimensional problems.

**Necmiye Ozay** received the B.S. degree from Bogazici University, Istanbul in 2004, the M.S. degree from the Pennsylvania State University, University Park in 2006 and the Ph.D. degree from Northeastern University, Boston in 2010, all in electrical engineering. She was a postdoctoral scholar at California Institute of Technology, Pasadena between 2010 and 2013. She is currently an associate professor of Electrical Engineering and Computer Science, at University of Michigan, Ann Arbor.