

Received XX September 2021; revised; accepted XX October 2021; Date of publication XX November 2021; date of current version XX November 2021.

Digital Object Identifier 10.1109/OJCSYS.2021.Doi Number

# State and Sparse Input Estimation in Linear Dynamical Systems using Low-Dimensional Measurements

RUPAM KALYAN CHAKRABORTY<sup>1</sup>, GEETHU JOSEPH<sup>1</sup> (Senior Member, IEEE), CHANDRA R. MURTHY<sup>2</sup> (Fellow, IEEE)

<sup>1</sup> Faculty of Electrical Engineering, Mathematics, and Computer Science, Delft University of Technology, Delft 2628 CD, The Netherlands; emails: {R.K.Chakraborty,G.Joseph}@tudelft.nl

ABSTRACT Sparsity constraints on the control inputs of a linear dynamical system naturally arise in several practical applications such as networked control, computer vision, seismic signal processing, and cyber-physical systems. In this work, we consider the problem of jointly estimating the states and sparse inputs of such systems from low-dimensional (compressive) measurements. Due to the low-dimensional measurements, conventional Kalman filtering and smoothing algorithms fail to accurately estimate the states and inputs. We present a Bayesian approach that exploits the input sparsity to significantly improve estimation accuracy. Sparsity in the input estimates is promoted by using different prior distributions on the input. We investigate two main approaches: regularizer-based maximum a posteriori estimation and Bayesian learning-based estimation. We also extend the approaches to handle control inputs with common support and analyze the time and memory complexities of the presented algorithms. Finally, using numerical simulations, we show that our algorithms outperform the state-of-the-art methods in terms of accuracy and time/memory complexities, especially in the low-dimensional measurement regime.

INDEX TERMS Kalman filtering and smoothing, sparsity-promoting regularizer, joint sparsity, ADMM,  $\ell_1$  minimization, reweighted  $\ell_2$  minimization, sparse Bayesian learning, variational Bayesian methods.

## I. Introduction

Sparse actuator control of linear dynamical systems (LDSs) has recently gained considerable interest in the literature [2]– [5]. This new research area deals with LDSs with sparsity constraints on the control inputs, i.e., each input vector contains only a small number of nonzero entries. An LDS with sparse control inputs models several practical applications such as networked control systems [2], [3], opinion dynamics manipulation [6], [7], computer vision [8], seismic signal processing [9], [10], and cyber-physical systems [11], [12]. In such applications, an important goal is to jointly estimate the states and sparse inputs of the LDS from its measurements or output. For example, malicious attacks on cyberphysical systems can be modeled as sparse inputs [11], [12]. Recovery of these sparse attacks is crucial to detecting and mitigating the attacks, as demonstrated in applications such as aircraft engine [13]–[15] and power system network [16].

Motivated by the above applications, this paper focuses on developing *joint state and input recovery algorithms* for observable LDSs with sparse control inputs. Specifically, we consider a discrete-time LDS, with state transition matrix  $A_k \in \mathbb{R}^{n \times n}$ , input matrix  $B_k \in \mathbb{R}^{n \times m}$ , measurement matrices  $C_k \in \mathbb{R}^{p \times n}$  and  $D_k \in \mathbb{R}^{p \times m}$  at time instant k, whose dynamics are governed by

$$\boldsymbol{x}_{k+1} = \boldsymbol{A}_k \boldsymbol{x}_k + \boldsymbol{B}_k \boldsymbol{u}_k + \boldsymbol{w}_k \tag{1}$$

$$\boldsymbol{y}_k = \boldsymbol{C}_k \boldsymbol{x}_k + \boldsymbol{D}_k \boldsymbol{u}_k + \boldsymbol{v}_k. \tag{2}$$

Here,  $u_k \in \mathbb{R}^m$  is the input,  $x_k \in \mathbb{R}^n$  is the state, and  $y_k \in \mathbb{R}^p$  is the measurement at time k. Also,  $w_k$  and  $v_k$  are the process noise and measurement noise, respectively. The noise  $w_k \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_k)$  and  $v_k \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_k)$  are independent, where  $\mathbf{Q}_k \in \mathbb{R}^{n \times n}$  and  $\mathbf{R}_k \in \mathbb{R}^{p \times p}$  are positive definite matrices. Here,  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  denotes the Gaussian distribution with mean vector  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$ .

This work is licensed under a Creative Commons Attribution 4.0 License. For more information, see https://creativecommons.org/licenses/by/4.0/

<sup>&</sup>lt;sup>2</sup>Department of Electrical Communication Engineering, Indian Institute of Science (IISc), Bangalore 560012, India; email: cmurthy@iisc.ac.in R.K. Chakraborty's work was initiated while at IISc, Bangalore. This work was presented in part in [1].

We aim to simultaneously estimate the states and sparse inputs  $\{\boldsymbol{x}_k, \boldsymbol{u}_k : \|\boldsymbol{u}_k\|_0 \ll m\}_{k=1}^K$  from the low dimensional measurements  $\{\boldsymbol{y}_k\}_{k=1}^K$  with p < m, for a given K > 0. Here,  $\|\cdot\|_0$  denotes the  $\ell_0$  norm. We emphasize that our focus is not on the design of sparse control inputs, but on recovering the states and sparse inputs of the system.

#### A. Related Work

Joint recovery of states and input without assuming any specific structure on the inputs or states has been studied extensively [17]–[21], but these works ignore any underlying sparsity structure that may exist in the system. Exploiting sparsity can facilitate the recovery of states and inputs with far fewer measurements than conventional approaches. The existing studies on the recovery of states and inputs in sparsity-constrained LDS consider one of three problems:

#### 1) Estimation of sparse initial state $x_1$

Estimating the sparse initial state of an LDS is equivalent to the standard sparse recovery problem, and can be solved using algorithms like LASSO, orthogonal matching pursuit, or sparse Bayesian learning (SBL) [22]–[24]. Theoretical results on the recoverability of a sparse initial state from low dimensional linear measurements of the states when the sparsity is exploited have been derived in [23]–[27]. These works focus on estimating the sparse initial state, assuming complete knowledge of the inputs applied to the system.

# 2) Estimation of sparse states $\{x_k\}_{k=1}^N$

The sparse state estimation of an LDS without the knowledge of inputs has been studied in diverse applications. Sparsityaware Kalman filtering was proposed in [28] to track abrupt changes in the sequence of sparse states of an LDS. Further, sparse state estimation was discussed in [29], where an SBL-based algorithm was used under the assumption of jointly sparse states. Another approach for jointly sparse state recovery algorithms imposed an  $\ell_1$ -regularizer in the Kalman smoothing cost, followed by the alternating direction method of multipliers (ADMM) method [30]. Furthermore, the recovery of sparse states without the joint sparsity assumption was studied via  $\ell_1$ -regularization-based dynamic filtering [31] and via a variational form of SBL in [32]. Additionally, [33] considered a general non-linear state space model with linear measurements of sparse states and developed an SBL-based dynamic filtering algorithm. However, these works assume the system matrices and inputs are restricted to ensure that  $x_{k+1}$  in (1) is sparse.

3) Joint estimation of states and sparse inputs  $\{x_k, u_k\}_{k=1}^N$  In [34], the problem of jointly recovering the state and sparse input sequences as an  $\ell_1$ -minimization using convex optimization methods was considered, i.e., minimizing  $\sum_{k=1}^K \|u_k\|_1$ . Necessary and sufficient conditions for observability of sparse control inputs and the initial state for a noiseless LDS have been investigated in [35], [36]. These methods involve solving for a large-dimensional unknown sparse vector obtained by stacking the input vectors

and do not exploit the temporal correlation in the state evolution. Consequently, the resulting algorithms have high computational complexity and memory requirements. To address these literature gaps, we presents new sparsity-driven estimators with better performance and lower complexities.

#### B. Contributions

We address the joint estimation of states and sparse inputs from (1) and (2) using p observations per time step over K time steps. With K(n+m) unknowns and only Kp observations with  $p \ll n+m$ , the system is highly underdetermined. We solve the problem by enforcing sparsity via fictitious priors from the exponential family and integrating the resulting Bayesian framework with Kalman smoothing to exploit temporal correlations, yielding different estimation methods. Our contributions are as follows:

- 1) Regularizer-based approach: We integrate the sparsity-promoting priors in the maximum a posteriori (MAP) estimator of the system's state and inputs leading to  $\ell_1$ -regularized and reweighted  $\ell_2$ -regularized algorithms. (See Section III.)
- 2) Bayesian learning-based approach: We develop two techniques that use hierarchical Gaussian priors to induce sparsity. First, we rely on type-II maximum likelihood (ML) estimation combined with Kalman smoothing to determine the states and sparse inputs. Second, we develop a variational Bayesian (VB) algorithm that groups the prior parameters and unknown states and inputs as unobserved variables, and infers their posterior distributions. (See Section IV.)
- 3) Comparison and extension: We analyze and derive the time and memory complexities of both approaches in Section IV.C. Our empirical studies in Section V show that Bayesian learning-based algorithms outperform the regularizer-based approach and have lower complexity. We also extend the two approaches to the case of jointly sparse control inputs and present a similar analysis.

The key innovation in this work is seamlessly integrating advanced sparse recovery techniques with the Kalman smoothing framework. Unlike most control algorithms that exploit sparsity using basic  $\ell_1$ -norm regularization, often neglecting the temporal evolution of the state process, our work bridges the gap between sparse signal processing and control theory by introducing sophisticated Bayesian sparse signal recovery methods within the Kalman smoothing framework.

This work significantly extends the conference paper [1] in several key aspects: we present a family of regularizer-based approaches, convergence analysis of the algorithms, and more comprehensive simulation results along with comparisons to extended versions of state-of-the-art algorithms.

**Notation:** The ith entry of vector  $\boldsymbol{a}$  is  $\boldsymbol{a}(i)$  and the (i,j)th entry of matrix  $\boldsymbol{A}$  is  $\boldsymbol{A}(i,j)$ . We denote a sequence of vectors  $\{\boldsymbol{a}_1,\boldsymbol{a}_2,\ldots,\boldsymbol{a}_N\}$  by  $\boldsymbol{A}_1^N$ . The  $\ell_1$ -norm is denoted by  $\|\cdot\|_1$ , and  $\|\cdot\|$  denotes the induced  $\ell_2$ -norm for matrices and Euclidean norm for vectors. Also,  $\boldsymbol{A}_{\mathcal{S}}$  denotes the submatrix of  $\boldsymbol{A}$  with the columns indexed by set  $\mathcal{S}$ , and  $\boldsymbol{x}_{\mathcal{S}}$  denotes the subvector of  $\boldsymbol{x}$  with entries indexed by  $\mathcal{S}$ . If  $\boldsymbol{P}$  is a nonsingular matrix, we define  $\|\boldsymbol{x}\|_{\boldsymbol{P}} = \boldsymbol{x}^T \boldsymbol{P}^{-1} \boldsymbol{x}$ .

#### II. MAP Estimation of States and Sparse Inputs

In this section, we introduce the generalized framework for Bayesian estimation of initial state and sparse inputs via sparsity-promoting priors. Specifically, we consider the estimation of the states and inputs  $\{x_k, u_k\}_{k=1}^K$  using measurements  $Y_1^K \triangleq \{y_k\}_{k=1}^K$  in (2) using a Bayesian framework, exploiting the sparsity of the inputs and the temporal correlation across the states dictated by (1).

The MAP estimates  $\hat{\boldsymbol{x}}_{k|K}$ ,  $\hat{\boldsymbol{u}}_{k|K}$  of the states and inputs are computed using the joint distribution of all the states and inputs  $\{\boldsymbol{x}_k, \boldsymbol{u}_k\}_{k=1}^K$  given all observations  $\boldsymbol{Y}_1^K$ , subject to the system dynamics in (1) and (2), i.e.,

$$\begin{aligned}
\left\{\hat{\boldsymbol{x}}_{k|K}, \hat{\boldsymbol{u}}_{k|K}\right\}_{k=1}^{K} &= \underset{\boldsymbol{x}_{k}, \boldsymbol{u}_{k}}{\operatorname{arg max}} p\left(\left\{\boldsymbol{x}_{k}, \boldsymbol{u}_{k}\right\}_{k=1}^{K} \mid \boldsymbol{Y}_{1}^{K}\right) \\
&= \underset{\boldsymbol{x}_{k}, \boldsymbol{u}_{k}}{\operatorname{arg max}} \prod_{k=1}^{K} p(\boldsymbol{y}_{k} | \boldsymbol{x}_{k}, \boldsymbol{u}_{k}) p(\boldsymbol{u}_{k}) \\
&\times \prod_{k=2}^{K} p(\boldsymbol{x}_{k} | \boldsymbol{x}_{k-1}, \boldsymbol{u}_{k-1}) \times p(\boldsymbol{x}_{1}), \quad (3)
\end{aligned}$$

where we use the Markov structure implied by (1), and assume the input  $\boldsymbol{u}_k$  is independent of the inputs. We note that the control inputs applied to the LDS need not be independent of the initial state. However, we make no assumptions about any prior knowledge of such relationships. Instead, we treat the inputs as unknown parameters and impose a fictitious prior to enforce sparsity (which we discuss later) to develop our algorithms.

Further, we do not make any assumptions about the initial state  $x_1$  and assume  $p(x_1)$  to be a (an improper) uniform prior. So, from (1) and (2), and using Gaussianity of the process and measurement noise, the optimization problem in (3) reduces to

$$\left\{\hat{\boldsymbol{x}}_{k|K}, \hat{\boldsymbol{u}}_{k|K}\right\}_{k=1}^{K} = \underset{k=1,...,K}{\operatorname{arg\,min}} \frac{1}{2} \sum_{k=1}^{K} \left\|\boldsymbol{y}_{k} - \boldsymbol{C}_{k} \boldsymbol{x}_{k} - \boldsymbol{D}_{k} \boldsymbol{u}_{k}\right\|_{\boldsymbol{R}_{k}}^{2}$$

$$+\frac{1}{2}\sum_{k=1}^{K-1} \|\boldsymbol{x}_{k+1} - \boldsymbol{A}_{k}\boldsymbol{x}_{k} - \boldsymbol{B}_{k}\boldsymbol{u}_{k}\|_{\boldsymbol{Q}_{k}}^{2} - \sum_{k=1}^{K} \ln(p(\boldsymbol{u}_{k})). \quad (4)$$

Since the input is known to be sparse, we encode this information into the estimation model via suitable priors on the inputs. Based on different sparsity-promoting priors, we develop two approaches: (a) regularized robust Kalman smoothing (RKS) and (b) Bayesian RKS, presented next.

# III. Regularized Robust Kalman Smoothing

Inspired by the convex optimization-based sparse signal recovery algorithms [22], [37], we use the following prior to induce sparsity:

$$p(\boldsymbol{u}_k) = \prod_{i=1}^{m} \chi \exp\left[-\frac{\tau_k}{2} |\boldsymbol{u}_k(i)|^l\right],$$

where  $\chi$  is the normalizing constant and  $\tau_k, l > 0$  are known distribution parameters. Here,  $\tau_k$  controls the sparsity of

the control inputs, i.e., a large value of  $\tau_k$  leads to sparser solutions. This parameter is often chosen by cross-validation. The choice of l determines the properties of the optimization problem in (4), leading to different estimators as given below.

## A. $\ell_1$ -regularized Robust Kalman Smoothing

Motivated by the  $\ell_1$  norm-based Laplacian prior, the most popular choice is l=1. With l=1, the optimization problem in (4) becomes

$$\left\{\hat{\boldsymbol{x}}_{k|K}, \hat{\boldsymbol{u}}_{k|K}\right\}_{k=1}^{K} = \underset{k=1, \dots, K}{\arg\min} \sum_{k=1}^{K} \left\|\boldsymbol{y}_{k} - \boldsymbol{C}_{k} \boldsymbol{x}_{k} - \boldsymbol{D}_{k} \boldsymbol{u}_{k}\right\|_{\boldsymbol{R}_{k}}^{2}$$

$$+ \sum_{k=1}^{K-1} \|\boldsymbol{x}_{k+1} - \boldsymbol{A}_{k} \boldsymbol{x}_{k} - \boldsymbol{B}_{k} \boldsymbol{u}_{k}\|_{\boldsymbol{Q}_{k}}^{2} + \sum_{k=1}^{K} \tau_{k} \|\boldsymbol{u}_{k}\|_{1}. \quad (5)$$

The convex optimization problem in (5) has no closed-form solution. So, we use the ADMM algorithm [38] to solve it. ADMM decomposes the convex optimization problem in (5) into simpler optimization problems. We reformulate (5) to another equivalent optimization problem using auxiliary variables  $\{t_k \in \mathbb{R}^m\}_{k=1}^K$ . We define the augmented Lagrangian function as

$$\mathcal{L}\left(\left\{\boldsymbol{x}_{k}, \boldsymbol{u}_{k}, \boldsymbol{t}_{k}, \boldsymbol{\lambda}_{k}\right\}_{k=1}^{K}\right) = \sum_{k=1}^{K} \left\|\boldsymbol{y}_{k} - \boldsymbol{C}_{k} \boldsymbol{x}_{k} - \boldsymbol{D}_{k} \boldsymbol{u}_{k}\right\|_{\boldsymbol{R}_{k}}^{2}$$

$$+ \sum_{k=1}^{K-1} \left\|\boldsymbol{x}_{k+1} - \boldsymbol{A}_{k} \boldsymbol{x}_{k} - \boldsymbol{B}_{k} \boldsymbol{u}_{k}\right\|_{\boldsymbol{Q}_{k}}^{2}$$

$$+ \sum_{k=1}^{K} \left[\tau_{k} \left\|\boldsymbol{t}_{k}\right\|_{1} + \boldsymbol{\lambda}_{k}^{\mathsf{T}} \left(\boldsymbol{t}_{k} - \boldsymbol{u}_{k}\right) + c \left\|\boldsymbol{t}_{k} - \boldsymbol{u}_{k}\right\|^{2}\right], \quad (6)$$

where  $\{\lambda_k \in \mathbb{R}^m\}_{k=1}^K$  are the Lagrangian multipliers that arise from the K constraints  $\boldsymbol{t}_k = \boldsymbol{u}_k$ . Also, c > 0 is a scalar. ADMM is an iterative algorithm that alternately solves subproblems of (6), each focusing on a specific block of variables. The rth iteration updates are

$$\left\{\boldsymbol{x}_{k}^{(r)}, \boldsymbol{u}_{k}^{(r)}\right\}_{k=1}^{K} = \underset{\left\{\boldsymbol{x}_{k}, \boldsymbol{u}_{k}\right\}_{k=1}^{K}}{\arg\min} \mathcal{L}\left(\left\{\boldsymbol{x}_{k}, \boldsymbol{u}_{k}, \boldsymbol{t}_{k}^{(r-1)} \boldsymbol{\lambda}_{k}^{(r-1)}\right\}_{k=1}^{K}\right)$$
(7

$$\boldsymbol{t}_{k}^{(r)} = \arg\min_{\boldsymbol{t}_{k}} \mathcal{L}\left(\left\{\boldsymbol{x}_{k}^{(r)}, \boldsymbol{u}_{k}^{(r)}, \boldsymbol{t}_{k}, \boldsymbol{\lambda}_{k}^{(r-1)}\right\}_{k=1}^{K}\right) \tag{8}$$

$$\boldsymbol{\lambda}_{k}^{(r)} = \boldsymbol{\lambda}_{k}^{(r-1)} + 2c\left(\boldsymbol{u}_{k}^{(r)} - \boldsymbol{t}_{k}^{(r)}\right), \tag{9}$$

for k = 1, 2, ..., K. Further, (7) can be simplified as

$$\begin{aligned} \left\{ \boldsymbol{x}_{k}^{(r)}, \boldsymbol{u}_{k}^{(r)} \right\}_{k=1}^{K} &= \underset{\boldsymbol{x}_{k}, \boldsymbol{u}_{k}}{\operatorname{arg \, min}} \sum_{k=1}^{K} \left\| \boldsymbol{y}_{k}^{\ell_{1}} - \tilde{\boldsymbol{C}}_{k} \boldsymbol{x}_{k} - \tilde{\boldsymbol{D}}_{k} \boldsymbol{u}_{k} \right\|_{\boldsymbol{R}_{k}^{\ell_{1}}}^{2} \\ &+ \sum_{k=1}^{K-1} \left\| \boldsymbol{x}_{k+1} - \boldsymbol{A}_{k} \boldsymbol{x}_{k} - \boldsymbol{B}_{k} \boldsymbol{u}_{k} \right\|_{\boldsymbol{Q}_{k}}^{2}. \quad (10)
\end{aligned}$$

Here, we define the new matrices

$$\mathbf{y}_{k}^{\ell_{1}} = \begin{bmatrix} \mathbf{y}_{k} \\ \mathbf{t}_{k}^{(r-1)} - c^{-1} \boldsymbol{\lambda}_{k}^{(r-1)} \end{bmatrix}, \quad \mathbf{R}_{k}^{\ell_{1}} = \begin{bmatrix} \mathbf{R}_{k} & \mathbf{0} \\ \mathbf{0} & c^{-1} \mathbf{I} \end{bmatrix}$$
(11)
$$\tilde{\mathbf{C}}_{k} = \begin{bmatrix} \mathbf{C}_{k}^{\mathsf{T}} & \mathbf{0} \end{bmatrix}^{\mathsf{T}}, \qquad \tilde{\mathbf{D}}_{k} = \begin{bmatrix} \mathbf{D}_{k}^{\mathsf{T}} & \mathbf{I} \end{bmatrix}^{\mathsf{T}}.$$
(12)

To solve (10), we next extend Kalman filtering and smoothing, laying the foundation for our algorithms in this paper.

#### Theorem 1:

Consider the MAP estimates of the states and inputs of the linear system in (1) and (2) with the assumption that  $D_k$ 's have full column rank and without any prior information on the inputs, given by

$$\underset{k=1,...,K}{\operatorname{arg \, min}} \sum_{k=1}^{K} \|\boldsymbol{y}_{k} - \boldsymbol{C}_{k}\boldsymbol{x}_{k} - \boldsymbol{D}_{k}\boldsymbol{u}_{k}\|_{\boldsymbol{R}_{k}}^{2} + \sum_{k=1}^{K-1} \|\boldsymbol{x}_{k+1} - \boldsymbol{A}_{k}\boldsymbol{x}_{k} - \boldsymbol{B}_{k}\boldsymbol{u}_{k}\|_{\boldsymbol{Q}_{k}}^{2}. \quad (13)$$

This problem can be solved recursively using the prediction, filtering, and smoothing steps, using the robust Kalman smoothing (RKS) algorithm summarized in Algorithm 1.

*Proof:* See Appendix A.

We note that the statement of Theorem 1 requires  $p \ge m$ , which implies that, in the absence of additional structural information about the inputs  $u_k$ , a necessary condition for obtaining MAP estimates is that the number of measurements must be at least as large as the input dimension.

In Algorithm 1, we determine the statistics of the joint Gaussian posterior distribution of  $x_k$  and  $u_k$  given  $Y_1^k$ . We know for any t and k the posterior distribution of the state  $x_t$  and input  $u_t$  given measurements  $Y_1^K$  is given by

$$p\left(\boldsymbol{x}_{t}, \boldsymbol{u}_{t} | \boldsymbol{Y}_{1}^{k}\right) = \mathcal{N}\left(\begin{bmatrix} \hat{\boldsymbol{x}}_{t|k} \\ \hat{\boldsymbol{u}}_{t|k} \end{bmatrix}, \begin{bmatrix} \boldsymbol{P}_{t|k}^{\boldsymbol{x}} & \boldsymbol{P}_{t|k}^{\boldsymbol{x}\boldsymbol{u}} \\ \left(\boldsymbol{P}_{t|k}^{\boldsymbol{x}\boldsymbol{u}}\right)^{\mathsf{T}} & \boldsymbol{P}_{t|k}^{\boldsymbol{u}} \end{bmatrix}\right). \quad (14)$$

Here,  $\hat{x}_{t|k} \in \mathbb{R}^n$  and  $\hat{u}_{t|k} \in \mathbb{R}^m$  are the estimates of  $x_t$  and  $u_t$  given  $Y_1^k$  with associated covariances  $P_{t|k}^x \in \mathbb{R}^{n \times n}$  and  $P_{t|k}^u \in \mathbb{R}^{m \times m}$ , respectively, and cross-covariance  $P_{t|k}^{xu} \in \mathbb{R}^{n \times m}$ . Then, defining  $\boldsymbol{\xi}_t^\mathsf{T} = \begin{bmatrix} \boldsymbol{x}_t^\mathsf{T} & \boldsymbol{u}_t^\mathsf{T} \end{bmatrix}$ , we get  $p\left(\boldsymbol{\xi}_t|Y_1^k\right) \sim \mathcal{N}\left(\hat{\boldsymbol{\xi}}_{t|k}, P_{t|k}^{\boldsymbol{\xi}}\right)$ , where  $\hat{\boldsymbol{\xi}}_{t|k}$  and  $P_{t|k}^{\boldsymbol{\xi}}$  are the mean and covariance of the distribution in (14).

Since (10) has the same form as (13), its solution follows the RKS algorithm with the measurement model in (11) and (12). Further, from (8), the auxiliary variable update is

$$\mathbf{t}_{k}^{(r)} = \arg\min_{\mathbf{t}_{k}} \ \tau_{k} \|\mathbf{t}_{k}\|_{1} + c \|\mathbf{t}_{k} - c^{-1}\boldsymbol{\lambda}_{k}^{(r-1)} - \boldsymbol{u}_{k}^{(r)}\|^{2}$$
$$= S_{c^{-1}\tau_{k}} \left(\boldsymbol{u}_{k}^{(r)} + c^{-1}\boldsymbol{\lambda}_{k}^{(r-1)}\right), \tag{15}$$

where the last step follows from the LASSO solution. Also,  $S(\cdot)$  is the entry-wise soft thresholding function,  $S_b(a) = \operatorname{Sgn}(a) \max\{|a| - b, 0\}$ . Combing (9), (10) with Theorem 1, and (15),  $\ell_1$ -regularized RKS is summarized in Algorithm 2.

## Algorithm 1 Robust Kalman Smoothing

Inputs: 
$$\{\boldsymbol{y}_k, \boldsymbol{A}_k, \boldsymbol{B}_k, \boldsymbol{C}_k, \boldsymbol{D}_k, \boldsymbol{Q}_k, \boldsymbol{R}_k\}_{k=1}^K$$
  
Initialization:  $\hat{\boldsymbol{\xi}}_{0|0} = \boldsymbol{0} \in \mathbb{R}^{n+m}, \, \boldsymbol{P}_{0|0}^{\boldsymbol{\xi}} \in \mathbb{R}^{(n+m)\times(n+m)}$   
and  $\boldsymbol{Q}_0 = \boldsymbol{0} \in \mathbb{R}^{n\times n}$   
1:  $\boldsymbol{\xi}_k = \begin{bmatrix} \boldsymbol{x}_k^\mathsf{T} & \boldsymbol{u}_k^\mathsf{T} \end{bmatrix}^\mathsf{T}, \, \tilde{\boldsymbol{A}}_k = \begin{bmatrix} \boldsymbol{A}_k & \boldsymbol{B}_k \end{bmatrix}$  and  $\boldsymbol{T} = \begin{bmatrix} \boldsymbol{I} & \boldsymbol{0} \end{bmatrix}^\mathsf{T} \in \mathbb{R}^{(n+m)\times n}$   
2:  $\boldsymbol{J}_k = \begin{pmatrix} \boldsymbol{D}_k^\mathsf{T} \boldsymbol{R}_k^{-1} \boldsymbol{D}_k \end{pmatrix}^{-1} \boldsymbol{D}_k^\mathsf{T} \boldsymbol{R}_k^{-1} \text{ for } k = 1, \dots K$ 

3: **for** k = 1, 2, ..., K **do** 

#Prediction:

4: 
$$\hat{x}_{k|k-1} = \tilde{A}_{k-1}\hat{\xi}_{k-1|k-1}$$
  
5:  $P_{k|k-1}^{x} = \tilde{A}_{k-1}P_{k-1|k-1}^{\xi}\tilde{A}_{k-1}^{\mathsf{T}} + Q_{k-1}$   
#Filtering:

6: 
$$L_{k} = \boldsymbol{P}_{k|k-1}^{\boldsymbol{x}} \boldsymbol{C}_{k}^{\mathsf{T}} \left( \boldsymbol{R}_{k} + \boldsymbol{C}_{k} \boldsymbol{P}_{k|k-1}^{\boldsymbol{x}} \boldsymbol{C}_{k}^{\mathsf{T}} \right)^{-1}$$
7: 
$$\boldsymbol{\mathcal{G}}_{k} = \begin{bmatrix} (\boldsymbol{I} - \boldsymbol{L}_{k} \boldsymbol{D}_{k} \boldsymbol{J}_{k} \boldsymbol{C}_{k})^{-1} \boldsymbol{L}_{k} (\boldsymbol{I} - \boldsymbol{D}_{k} \boldsymbol{J}_{k}) \\ (\boldsymbol{I} - \boldsymbol{J}_{k} \boldsymbol{C}_{k} \boldsymbol{L}_{k} \boldsymbol{D}_{k})^{-1} \boldsymbol{J}_{k} (\boldsymbol{I} - \boldsymbol{C}_{k} \boldsymbol{L}_{k}) \end{bmatrix}$$
8: 
$$\hat{\boldsymbol{\xi}}_{k|k} = \boldsymbol{T} \hat{\boldsymbol{x}}_{k|k-1} + \boldsymbol{\mathcal{G}}_{k} \left( \boldsymbol{y}_{k} - \boldsymbol{C}_{k} \hat{\boldsymbol{x}}_{k|k-1} \right)$$
9: 
$$\boldsymbol{P}_{k|k}^{\boldsymbol{\xi}} = (\boldsymbol{T} - \boldsymbol{\mathcal{G}}_{k} \boldsymbol{C}_{k}) \boldsymbol{P}_{k|k-1}^{\boldsymbol{x}} (\boldsymbol{T} - \boldsymbol{\mathcal{G}}_{k} \boldsymbol{C}_{k})^{\mathsf{T}} + \boldsymbol{\mathcal{G}}_{k} \boldsymbol{R}_{k} \boldsymbol{\mathcal{G}}_{k}^{\mathsf{T}}$$

Compute  $\hat{x}_{k|k}$  and  $P_{k|k}^{x}$  from  $\hat{\xi}_{k|k}$  and  $P_{k|k}^{\xi}$ 

11: **end for** 

17: **end for** 

10:

#Smoothing:

12: **for** 
$$k = K - 1, K - 2, \dots, 1$$
 **do**
13:  $\boldsymbol{K}_{k} = \boldsymbol{P}_{k|k}^{\boldsymbol{\xi}} \tilde{\boldsymbol{A}}_{k}^{\mathsf{T}} \left( \boldsymbol{P}_{k+1|k}^{\boldsymbol{x}} \right)^{-1}$ 
14:  $\boldsymbol{P}_{k|K}^{\boldsymbol{\xi}} = \boldsymbol{P}_{k|k}^{\boldsymbol{\xi}} + \boldsymbol{K}_{k} \left( \boldsymbol{P}_{k+1|K}^{\boldsymbol{x}} - \boldsymbol{P}_{k+1|k}^{\boldsymbol{x}} \right) \boldsymbol{K}_{k}^{\mathsf{T}}$ 
15:  $\hat{\boldsymbol{\xi}}_{k|K} = \hat{\boldsymbol{\xi}}_{k|k} + \boldsymbol{K}_{k} \left( \hat{\boldsymbol{x}}_{k+1|K} - \tilde{\boldsymbol{A}}_{k} \hat{\boldsymbol{\xi}}_{k|k} \right)$ 
16: Compute  $\hat{\boldsymbol{x}}_{k|K}$  and  $\boldsymbol{P}_{k|K}^{\boldsymbol{x}}$  from  $\hat{\boldsymbol{\xi}}_{k|K}$  and  $\boldsymbol{P}_{k|K}^{\boldsymbol{\xi}}$ 

Outputs:  $\{\hat{\pmb{x}}_{k|K}\}_{k=1}^K$  and  $\{\hat{\pmb{u}}_{k|K}\}_{k=1}^K$ 

In Algorithm 2, the full column rank requirement of Theorem 1 holds for  $\tilde{\boldsymbol{D}}_k = \begin{bmatrix} \boldsymbol{D}_k^\mathsf{T} & \boldsymbol{I} \end{bmatrix}^\mathsf{T} \in \mathbb{R}^{(p+m) \times m}$ .

Since (5) is convex, ADMM is guaranteed to converge to the global solution. A general limitation of  $\ell_1$ -based regularization, however, is that it may yield biased estimates for large coefficients, thereby necessitating careful tuning of the penalty parameters [39].

## B. Reweighted $\ell_2$ -regularized Robust Kalman Smoothing

When 0 < l < 2, an alternative algorithm, similar in spirit to the iterative reweighting-based sparse recovery [37], can be derived. Then, the optimization problem (13) changes to

$$\left\{ \hat{\boldsymbol{x}}_{k|K}, \hat{\boldsymbol{u}}_{k|K} \right\}_{k=1}^{K} = \underset{k=1, \dots, K}{\arg\min} \sum_{k=1}^{K} \left\| \boldsymbol{y}_{k} - \boldsymbol{C}_{k} \boldsymbol{x}_{k} - \boldsymbol{D}_{k} \boldsymbol{u}_{k} \right\|_{\boldsymbol{R}_{k}}^{2}$$

$$+\sum_{k=1}^{K-1} \|\boldsymbol{x}_{k+1} - \boldsymbol{A}_{k} \boldsymbol{x}_{k} - \boldsymbol{B}_{k} \boldsymbol{u}_{k}\|_{\boldsymbol{Q}_{k}}^{2} + \sum_{k=1}^{K} \tau_{k} \sum_{i=1}^{m} |\boldsymbol{u}_{k}(i)|^{l}.$$
 (16)

# **Algorithm 2** $\ell_1$ -regularized Robust Kalman Smoothing

Inputs: 
$$\{y_k, A_k, B_k, C_k, D_k, Q_k, R_k\}_{k=1}^K$$

**Parameters:** c,  $r_{\text{max}}$ , and  $\tau_k$  for k = 1, 2, ..., K, r = 1,

 $\epsilon = 2\epsilon_{\rm thres}$ 

Initialization:  $\boldsymbol{t}_k^{(0)} = \boldsymbol{\lambda}_k^{(0)} = \boldsymbol{0}$ , for  $k = 1, 2, \dots, K$ 

1: Compute  $R_k^{\ell_1}$ ,  $\tilde{C}_k$ ,  $\tilde{D}_k$  using (11) and (12)

2: while  $(r < r_{max})$  or  $(\epsilon > \epsilon_{thres})$  do

Compute  $\boldsymbol{y}_k^{\ell_1}$  using (11)

Compute  $x_k^{(r)}$  and  $u_k^{(r)}$  via Algorithm 1 replacing  $m{y}_k, \, m{C}_k, \, m{D}_k, \, ext{and} \, m{R}_k \, ext{ with } m{y}_k^{\ell_1}, \, ilde{m{C}}_k, \, ilde{m{D}}_k, \, ext{and} \, m{R}_k^{\ell_1}$ 

5:

, respectively for 
$$k=1,2,\ldots,K$$
 do 
6: Using (15),  $\boldsymbol{t}_k^{(r)}=S_{c^{-1}\tau_k}\left(\boldsymbol{u}_k^{(r)}+c^{-1}\boldsymbol{\lambda}_k^{(r-1)}\right)$ 

7: 
$$\boldsymbol{\lambda}_k^{(r)} = \boldsymbol{\lambda}_k^{(r-1)} + 2c\left(\boldsymbol{u}_k^{(r)} - \boldsymbol{t}_k^{(r)}\right)$$

9: 
$$\epsilon = \sum_{k=1}^{K} \left\| \boldsymbol{\xi}_k^{(r)} - \boldsymbol{\xi}_k^{(r-1)} \right\|_2^2$$
 where  $\boldsymbol{\xi}_k^{(r)} = \begin{bmatrix} \boldsymbol{x}_k^{(r)} \\ \boldsymbol{u}_k^{(r)} \end{bmatrix}$ 

11: end while

Outputs: 
$$\left\{m{x}_k^{(r)}
ight\}_{k=1}^K$$
 and  $\left\{m{u}_k^{(r)}
ight\}_{k=1}^K$ 

We use the majorization minimization (MM) technique to solve (16). In MM, the cost function is upper-bounded around the current iterate by a tractable function that is easy to optimize, and we minimize that upper bound to get the new iterate. Here, we bound the regularizer term by noting that the function  $|u|^l = (|u|^2)^{l/2}$  is a concave function of  $|u|^2$ . Therefore, the function is bounded above by its firstorder Taylor approximation, and for any  $u \in \mathbb{R}$ ,

$$(|\boldsymbol{u}_k(i)|^2)^{l/2} \le (|u|^2)^{l/2} + \frac{l}{2}(|u|^2)^{l/2-1}(|\boldsymbol{u}_k(i)|^2 - |u|^2).$$

With  $u = \boldsymbol{u}_h^{(r-1)}(i)$ , we get

$$(|\boldsymbol{u}_{k}(i)|^{2})^{l/2} \leq \frac{l}{2} |\boldsymbol{u}_{k}^{(r-1)}(i)|^{-(2-l)} |\boldsymbol{u}_{k}(i)|^{2} + (1 - \frac{l}{2}) |\boldsymbol{u}_{k}^{(r-1)}(i)|^{l}. \quad (17)$$

We now iteratively solve (16) by optimizing the upper bound in each iteration. So, the rth iteration computes

$$+ \sum_{k=1}^{K-1} \|\boldsymbol{x}_{k+1} - \boldsymbol{A}_{k} \boldsymbol{x}_{k} - \boldsymbol{B}_{k} \boldsymbol{u}_{k}\|_{\boldsymbol{Q}_{k}}^{2} + \frac{l}{2} \sum_{k=1}^{K} \tau_{k} \|\boldsymbol{u}_{k}\|_{\boldsymbol{W}_{k}^{(r)}}^{2}.$$
(18)

Here, the diagonal weight matrix  $\boldsymbol{W}_{k}^{(r)} \in \mathbb{R}^{m \times m}$  is

$$\boldsymbol{W}_{k}^{(r)} = \operatorname{diag}\left\{\left|\boldsymbol{u}_{k}^{(r-1)}\right|\right\}^{(2-l)},$$

with  $\left|\cdot\right|$  representing the element-wise modulus. When an entry of  $oldsymbol{u}_k^{(r-1)}$  become close to zeros,  $oldsymbol{W}_k^{(r)}$  can potentially become non-invertible. To avoid numerical instabilities, we prune the entries of  $oldsymbol{u}_k^{(r-1)}$  falls below some small threshold  $\kappa$  (e.g.,  $\kappa = 10^{-6}$ ) in magnitude. Now, (18) modifies to

$$\begin{aligned}
\left\{ \boldsymbol{x}_{k}^{(r)}, \tilde{\boldsymbol{u}}_{k}^{(r)} \right\}_{k=1}^{K} &= \underset{\boldsymbol{x}_{k}, \tilde{\boldsymbol{u}}_{k}}{\operatorname{arg \, min}} \sum_{k=1}^{K} \left\| \boldsymbol{y}_{k}^{\ell_{2}} - \boldsymbol{C}_{k}^{\ell_{2}} \boldsymbol{x}_{k} - \boldsymbol{D}_{k}^{\ell_{2}} \tilde{\boldsymbol{u}}_{k} \right\|_{\boldsymbol{R}_{k}^{\ell_{2}}}^{2} \\
&+ \sum_{k=1}^{K-1} \left\| \boldsymbol{x}_{k+1} - \boldsymbol{A}_{k} \boldsymbol{x}_{k} - \boldsymbol{B}_{k}^{\ell_{2}} \tilde{\boldsymbol{u}}_{k} \right\|_{\boldsymbol{Q}_{k}}^{2}, \quad (19)
\end{aligned}$$

where we define

$$\mathbf{y}_{k}^{\ell_{2}} = \begin{bmatrix} \mathbf{y}_{k}^{\mathsf{T}} & \mathbf{0}_{|\mathcal{I}_{k}|} \end{bmatrix}^{\mathsf{T}}, \qquad \mathbf{B}_{k}^{\ell_{2}} = (\mathbf{B}_{k})_{\mathcal{I}_{k}} \qquad (20)$$

$$\mathbf{C}_{k}^{\ell_{2}} = \begin{bmatrix} \mathbf{C}_{k}^{\mathsf{T}} & \mathbf{0}_{n \times |\mathcal{I}_{k}|} \end{bmatrix}^{\mathsf{T}}, \qquad \mathbf{D}_{k}^{\ell_{2}} = \begin{bmatrix} (\mathbf{D}_{k})_{\mathcal{I}_{k}}^{\mathsf{T}} & \mathbf{I}_{|\mathcal{I}_{k}|} \end{bmatrix}^{\mathsf{T}}, \qquad (21)$$

where  $\tilde{\boldsymbol{u}}_k^{(r-1)} \in \mathbb{R}^{|\mathcal{I}_k|}$  is the vector obtained after removing the entries of  $\boldsymbol{u}_k^{(r-1)}$  whose absolute values are less than  $\kappa$ , and  $\mathcal{I}_k$  is the set of indices corresponding to the entries of  $\boldsymbol{u}_k^{(r-1)}$  which have absolute value at least  $\kappa$ . Also, we define

$$\boldsymbol{R}_{k}^{\ell_{2}} = \begin{bmatrix} \boldsymbol{R}_{k} & \boldsymbol{0} \\ \boldsymbol{0} & \frac{2}{\tau_{k}l} \operatorname{diag} \left\{ \left| \tilde{\boldsymbol{u}}_{k}^{(r-1)} \right| \right\}^{(2-l)} \end{bmatrix}.$$
 (22)

Here, (19) takes the same form as (5) and can be solved via RKS in Algorithm 1. Finally, we reconstruct  $u_k^{(r)}$  as

$$\left(\boldsymbol{u}_{k}^{(r)}\right)_{\mathcal{I}_{k}} = \tilde{\boldsymbol{u}}_{k}^{(r)} \quad \text{and} \quad \left(\boldsymbol{u}_{k}^{(r)}\right)_{\mathcal{I}_{k}^{0}} = \boldsymbol{0}.$$
 (23)

The pseudocode for the overall reweighted  $\ell_2$ -regularized RKS is summarized in Algorithm 3.

**Algorithm 3** Reweighted  $\ell_2$ -regularized Robust Kalman Smoothing

Inputs:  $\{y_k, A_k, B_k, C_k, D_k, Q_k, R_k\}_{k=1}^K$ 

**Parameters:**  $l, r_{\text{max}}, \kappa, \text{ and } \tau_k \text{ for } k = 1, 2, \dots, K, \epsilon = 0$ 

**Initialization:**  $u_k^{(0)} = 1$ , for  $k = 1, 2, \dots, K$ 

1: while  $(r < r_{max})$  or  $(\epsilon > \epsilon_{thres})$  do

Determine  $\mathcal{I}_k$  for each  $oldsymbol{u}_k^{(r-1)}$  and compute  $ilde{oldsymbol{u}}_k^{(r-1)}$ 

Compute  $y_k^{\ell_2}$ ,  $B_k^{\ell_2}$ ,  $C_k^{\ell_2}$ ,  $D_k^{\ell_2}$ ,  $R_k^{\ell_2}$  using (20), (21) and (22)

4: Compute  $x_k^{(r)}$  and  $\tilde{u}_k^{(r)}$  via Algorithm 1 replacing  $y_k$ ,  $B_k,\,C_k,\,D_k$ , and  $R_k$  with  $m{y}_k^{\ell_2},\,B_k^{\ell_2},\,C_k^{\ell_2},\,D_k^{\ell_2}$ , and  $oldsymbol{R}_k^{\ell_2}$  respectively

5: Construct  $u_k^{(r)}$  using (23)

6: 
$$\epsilon = \sum_{k=1}^{K} \left\| \boldsymbol{\xi}_k^{(r)} - \boldsymbol{\xi}_k^{(r-1)} \right\|_2^2 \text{ where } \boldsymbol{\xi}_k^{(r)} = \begin{bmatrix} \boldsymbol{x}_k^{(r)} \\ \boldsymbol{u}_k^{(r)} \end{bmatrix}$$

Outputs: 
$$\left\{m{x}_k^{(r)}
ight\}_{k=1}^K$$
 and  $\left\{m{u}_k^{(r)}
ight\}_{k=1}^K$ 

Next, we show that iteratively optimizing the surrogate function in (18) via Algorithm 3 either decreases the value of the cost function in (16) or leaves it unchanged, and that the sequence of cost function values over iterations converges.

# **Proposition 1:**

For a given set of inputs, let  $\left\{x_k^{(r)}, u_k^{(r)}\right\}_{k=1}^K$  be the sequence generated by Algorithm 3. Then, the corresponding sequence of cost function values in (5) is monotonically nonincreasing and converges to a limit in  $\mathbb{R}$ .

#### Proof:

For brevity, we define  $x = X_1^K$  and  $u = U_1^K$ . Let the cost function in (16) be denoted by  $f(x, u) = f_1(x, u) +$  $f_2(\boldsymbol{u})$ , where  $f_1(\boldsymbol{x}, \boldsymbol{u})$  denotes the first two quadratic terms and  $f_2(\boldsymbol{u}) = \sum_{k=1}^K \tau_k \sum_{i=1}^m |\boldsymbol{u}_k(i)|^l$  is the regularizer term. Also, let the upper bound in (17) be  $g(\boldsymbol{u}|\boldsymbol{u}^{(r-1)}) \geq f_2(\boldsymbol{u})$  and equality holds when  $\boldsymbol{u} = \boldsymbol{u}^{(r-1)}$ . Hence, we derive

$$f(\boldsymbol{x}^{(r-1)}, \boldsymbol{u}^{(r-1)}) = f_1(\boldsymbol{x}^{(r-1)}, \boldsymbol{u}^{(r-1)}) + g(\boldsymbol{u}^{(r-1)}|\boldsymbol{u}^{(r-1)})$$

$$\geq f_1(\boldsymbol{x}^{(r)}, \boldsymbol{u}^{(r)}) + g(\boldsymbol{u}^{(r)}|\boldsymbol{u}^{(r)})$$

$$\geq f_1(\boldsymbol{x}^{(r)}, \boldsymbol{u}^{(r)}) + f_2(\boldsymbol{u}^{(r)})$$

$$= f(\boldsymbol{x}^{(r)}, \boldsymbol{u}^{(r)}).$$

The first inequality holds because, in the rth iteration,  $(\boldsymbol{x}^{(r)}, \boldsymbol{u}^{(r)})$  the RKS algorithm minimizes the cost function in (18), which is equivalent to minimizing  $f_1(\boldsymbol{x}, \boldsymbol{u}) + g(\boldsymbol{u}|\boldsymbol{u}^{(r-1)})$ . The second inequality holds because  $g(\boldsymbol{u}|\boldsymbol{u}^{(r)}) \geq f_2(\boldsymbol{u})$  for any  $\boldsymbol{u}$ . Hence, the sequence of the cost function values is monotonically non-increasing. Finally, f(x, u) is lower bounded by 0. Thus, by the monotone convergence theorem, the cost function sequence converges.

The above result guarantees that cost function derived from the iterates converges to a limit, but the iterates may oscillate, cycle, or converge to a local minimum or a saddle point.

Iterative reweighted  $\ell_2$  regularization provides stronger sparsity promotion than  $\ell_1$ -based regularizer, since 0 < l < 1offers a closer approximation to the  $\ell_0$  norm than  $\ell_1$ . However, for 0 < l < 1 the problem becomes non-convex, and a key drawback is that the algorithm may get trapped in local minima, as indicated by Proposition 1.

#### C. Extension to Jointly Sparse Inputs

We can extend  $\ell_1$ -regularized RKS to estimate jointly sparse inputs by modifying (5) as follows:

$$\left\{\hat{m{x}}_{k|K}, \hat{m{u}}_{k|K}
ight\}_{k=1}^{K} = rgmin_{m{x}_{k}, m{u}_{k}} \sum_{k=1}^{K} \left\|m{y}_{k} - m{C}_{k}m{x}_{k} - m{D}_{k}m{u}_{k}
ight\|_{m{R}_{k}}^{2}$$

$$+\sum_{k=1}^{K-1} \|\boldsymbol{x}_{k+1} - \boldsymbol{A}_k \boldsymbol{x}_k - \boldsymbol{B}_k \boldsymbol{u}_k\|_{\boldsymbol{Q}_k}^2 + \tau \sum_{i=1}^m \sqrt{\sum_{k=1}^K \boldsymbol{u}_k^2(i)}, \quad (24)$$

where regularizer  $\tau \sum_{i=1}^{m} \sqrt{\sum_{k=1}^{K} u_k^2(i)}$  is inspired by the LASSO type regularization [40]. As  $\tau > 0$  increases, the common support output by the algorithm shrinks, making the inputs sparser. Now, using auxiliary variables  $t_k$ 's, we reformulate (24) and solve the following using ADMM:

$$\begin{aligned}
& \left\{ \hat{\boldsymbol{x}}_{k|K}, \hat{\boldsymbol{u}}_{k|K} \right\}_{k=1}^{K} = \underset{\boldsymbol{x}_{k}, \boldsymbol{u}_{k}}{\operatorname{arg \, min}} \sum_{k=1}^{K} \|\boldsymbol{y}_{k} - \boldsymbol{C}_{k} \boldsymbol{x}_{k} - \boldsymbol{D}_{k} \boldsymbol{u}_{k} \|_{\boldsymbol{R}_{k}}^{2} \\
&+ \sum_{k=1}^{K-1} \|\boldsymbol{x}_{k+1} - \boldsymbol{A}_{k} \boldsymbol{x}_{k} - \boldsymbol{B}_{k} \boldsymbol{u}_{k} \|_{\boldsymbol{Q}_{k}}^{2} + \tau \sum_{i=1}^{m} \sqrt{\sum_{k=1}^{K} t_{k}^{2}(i)} \\
&+ \sum_{k=1}^{K} \boldsymbol{\lambda}_{k}^{\mathsf{T}} \left( \boldsymbol{t}_{k} - \boldsymbol{u}_{k} \right) + c \sum_{k=1}^{K} \|\boldsymbol{t}_{k} - \boldsymbol{u}_{k} \|^{2}, \quad (25)
\end{aligned}$$

where  $\{ oldsymbol{\lambda}_k \in \mathbb{R}^m \}_{k=1}^K$  are the Lagrangian multipliers that arise from the equality constraints  $t_k = u_k$  and c > 0 is a positive scalar. Since (25) is identical to (6) except for the terms in  $t_k$ , our modified  $\ell_1$ -regularized RKS is identical to Algorithm 2 except for Step 6, which changes to [40],

$$\begin{aligned} \boldsymbol{t}^{(r)}(i,:) &= \frac{\boldsymbol{u}^{(r)}(i,:) + c^{-1} \boldsymbol{\lambda}^{(r-1)}(i,:)}{\left\| \boldsymbol{u}^{(r)}(i,:) + c^{-1} \boldsymbol{\lambda}^{(r-1)}(i,:) \right\|} \\ &\times S_{c^{-1}\tau} \left( \left\| \boldsymbol{u}^{(r)}(i,:) + c^{-1} \boldsymbol{\lambda}^{(r-1)}(i,:) \right\| \right), \end{aligned}$$

where function S is defined in (15), and  $t^{(r)}(i,:)$ ,  $u^{(r)}(i,:)$ ,  $\lambda^{(r-1)}(i,:)$  follow the definition.

$$a(i,:) = \begin{bmatrix} a_1(i) & a_2(i) & \dots & a_K(i) \end{bmatrix}^\mathsf{T} \in \mathbb{R}^K,$$
 (26) for any set  $\{a_k \in \mathbb{R}^m\}_{k=1}^K$  and  $i = 1, 2, \dots, K$ . This modified version is referred to as group  $\ell_1$ -regularized RKS.

#### IV. Bayesian Robust Kalman Smoothing

In this section, we present an alternative approach, called Bayesian RKS, to estimate the states and sparse inputs  $\{\boldsymbol{x}_k, \boldsymbol{u}_k\}_{k=1}^K$  using measurements  $\{\boldsymbol{y}_k\}_{k=1}^K$  in (2). Unlike the regularized RKS that assumes the knowledge of the parameters of the sparse input prior, the Bayesian approach uses hierarchical priors which account for the uncertainty in the prior distribution. Specifically, we use a hierarchical Gaussian prior on the inputs to promote sparsity:

$$p(\boldsymbol{u}_k; \boldsymbol{\gamma}_k) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi\gamma_k(i)}} \exp\left(-\frac{\boldsymbol{u}_k(i)^2}{2\boldsymbol{\gamma}_k(i)}\right), \quad (27)$$

where  $\gamma_k \in \mathbb{R}^m$  is the unknown hyperparameter of the distribution. The Bayesian RKS learns the hyperparameters from the measurements, unlike the regularized RKS. We present two Bayesian RKS variants: SBL-RKS and VB-RKS.

## A. Sparse Bayesian Learning-based RKS

In the SBL framework, we first compute the ML estimate 
$$\hat{\gamma}_{k}^{\text{ML}}$$
 of the hyperparameter as 
$$\{\hat{\gamma}_{k}^{\text{ML}}\}_{k=1}^{K} = \underset{\substack{\gamma_{k} \in \mathbb{R}_{+}^{m \times 1} \\ k=1}}{\operatorname{max}} p\left(\boldsymbol{Y}_{1}^{K}; \{\boldsymbol{\gamma}_{k}\}_{k=1}^{K}\right). \tag{28}$$

Using the estimate  $\hat{\gamma}_k^{\text{ML}}$ , we can estimate the states and inputs using the Kalman filtering and smoothing algorithm. For this, we note that (1) and (2) are equivalent to

$$\boldsymbol{\xi}_{k+1} = \begin{bmatrix} \tilde{\boldsymbol{A}}_k \\ \mathbf{0} \end{bmatrix} \boldsymbol{\xi}_k + \begin{bmatrix} \boldsymbol{w}_k \\ \boldsymbol{z}_k \end{bmatrix} \text{ and } \boldsymbol{y}_k = \begin{bmatrix} \boldsymbol{C}_k & \boldsymbol{D}_k \end{bmatrix} \boldsymbol{\xi}_k + \boldsymbol{v}_k,$$
(29)

where  $\xi_k$  is defined as

$$\boldsymbol{\xi}_k = \begin{bmatrix} \boldsymbol{x}_k^\mathsf{T} & \boldsymbol{u}_k^\mathsf{T} \end{bmatrix}^\mathsf{T},\tag{30}$$

and  $z_k = u_{k+1}$  is an auxiliary variable. From (27), we get

$$\begin{bmatrix} \boldsymbol{w}_k \\ \boldsymbol{z}_k \end{bmatrix} \sim \mathcal{N} \left( \boldsymbol{0}, \begin{bmatrix} \boldsymbol{Q}_k & \boldsymbol{0} \\ \boldsymbol{0} & \operatorname{Diag} \left\{ \boldsymbol{\gamma}_{k+1} \right\} \end{bmatrix} \right).$$

Given  $\{\hat{\gamma}_k^{\mathrm{ML}}\}_{k=1}^K$ , estimating the states and inputs reduces to estimating  $\{\boldsymbol{\xi}_k\}_{k=1}^K$  using  $\boldsymbol{Y}_1^K$  via Kalman filtering and smoothing, due to the Gaussian assumptions.

Since optimization problem in (28) lack a closed-form solution, we employ expectation-maximization (EM), an iterative method with the expectation (E) and maximization (M) steps. In the rth iteration, the E-step computes the expected log-likelihood function  $\mathcal{Q}^{(r)}$  of  $\{\gamma_k\}_{k=1}^K$  with respect to  $\{\xi_k\}_{k=1}^K$  given data  $Y_1^K$  and the current estimate  $\gamma_k^{(r-1)}$  of the hyperparameter  $\gamma_k$  obtained in the previous iteration. The M-step maximizes the expected log-likelihood to obtain the new estimate of  $\gamma_k$ . Further, from the state space model in (29), the distribution of the data is given by

$$p\left(\!\boldsymbol{Y}_{1}^{K}, \left\{\boldsymbol{\xi}_{k}\right\}_{k=1}^{K}; \left\{\boldsymbol{\gamma}_{k}\right\}_{k=1}^{K}\!\right) \!=\! \prod_{k=1}^{K} \! p(\boldsymbol{y}_{k}|\boldsymbol{\xi}_{k}) p\left(\boldsymbol{\xi}_{k}|\boldsymbol{\xi}_{k-1}; \boldsymbol{\gamma}_{k}\right),$$

where  $\xi_0 = 0$ . Thus, the E-step is given by

$$\begin{aligned} \mathcal{Q}^{(r)} \left( \left\{ \boldsymbol{\gamma}_{k} \right\}_{k=1}^{K} \mid \left\{ \boldsymbol{\gamma}_{k}^{(r-1)} \right\}_{k=1}^{K} \right) \\ &= \sum_{k=1}^{K} \mathbb{E}_{\boldsymbol{\xi}_{k}, \boldsymbol{\xi}_{k-1} \mid \boldsymbol{Y}_{1}^{K}; \boldsymbol{\gamma}_{k}^{(r-1)}} \left\{ \log p(\boldsymbol{y}_{k} | \boldsymbol{\xi}_{k}) \right. \\ &\qquad \times p\left( \boldsymbol{x}_{k} \mid \boldsymbol{\xi}_{k-1} \right) p\left( \boldsymbol{u}_{k}; \boldsymbol{\gamma}_{k} \right) \right\}. \end{aligned}$$

From the above relation, the M-step that maximizes  $\mathcal{Q}^{(r)}$  with respect to  $\{\gamma_k\}_{k=1}^K$  is separable, and ignoring the terms independent of  $\gamma_k$ , the M-step reduces to

$$\boldsymbol{\gamma}_{k}^{\left(r+1\right)} = \operatorname*{arg\,max}_{\boldsymbol{\gamma}_{k}} \mathbb{E}_{\boldsymbol{u}_{k} \mid \boldsymbol{Y}_{1}^{K}; \boldsymbol{\gamma}_{k}^{\left(r-1\right)}} \left\{ p\left(\boldsymbol{u}_{k}; \boldsymbol{\gamma}_{k}\right) \right\}.$$

Using (27), we derive the M-step as

$$\begin{split} \boldsymbol{\gamma}_k^{(r+1)} &= \mathop{\arg\min}_{\boldsymbol{\Gamma} = \operatorname{Diag}\{\boldsymbol{\gamma}_k\}} \log |\boldsymbol{\Gamma}| + \operatorname{Tr} \left\{ \boldsymbol{\Gamma}^{-1} \! (\hat{\boldsymbol{u}}_{k|K} \hat{\boldsymbol{u}}_{k|K}^\mathsf{T} + \boldsymbol{P}_{k|K}^{\boldsymbol{u}}) \right\} \\ &= \operatorname{Diag} \left\{ \hat{\boldsymbol{u}}_{k|K} \hat{\boldsymbol{u}}_{k|K}^\mathsf{T} + \boldsymbol{P}_{k|K}^{\boldsymbol{u}} \right\}. \end{split}$$

Further,  $\hat{u}_{k|K}$  and  $P^{u}_{k|K}$  can be computed by applying Kalman filtering and smoothing on the modified state space model in (29). The overall SBL-RKS algorithm is summarized in Algorithm 4.

Next, we establish a result similar to Proposition 1 for SBL-RKS. We show that the EM updates in Algorithm 4

are guaranteed to monotonically increase (or maintain) the data log-likelihood in each iteration. Hence, the sequence of objective function values over iterations converges.

# Algorithm 4 RKS with Sparse Bayesian Learning

Inputs: 
$$\{y_k, A_k, B_k, C_k, D_k, Q_k, R_k\}_{k=1}^K$$
  
Parameters:  $r_{\max}$ ,  $\epsilon_{\text{thres}}$   
Initialization:  $\gamma_k^{(0)} = 1$  for  $k = 1, 2, ..., K$ ,  $r = 1$ ,  $\epsilon = 2\epsilon_{\text{thres}}$   
1:  $\bar{A}_k = \begin{bmatrix} A_k & B_k \\ 0 & 0 \end{bmatrix} \in \mathbb{R}^{(n+m)\times(n+m)}$  and  $\bar{C}_k = \begin{bmatrix} C_k & D_k \end{bmatrix}$   
2: while  $(r < r_{\max})$  or  $(\epsilon > \epsilon_{\text{thres}})$  do #E-Step:  
3:  $\hat{\xi}_{0|0} = 0$ ,  $P_{0|0}^{\xi} = I$   
4: for  $k = 1, 2, ..., K$  do  
5:  $\bar{Q}_{k-1} = \begin{bmatrix} Q_{k-1} & 0 \\ 0 & \text{Diag} \{\gamma_k\} \end{bmatrix}$  #Prediction:  
6:  $\hat{\xi}_{k|k-1} = \bar{A}_{k-1}\hat{\xi}_{k-1|k-1}$   
7:  $P_{k|k-1}^{\xi} = \bar{A}_{k-1}P_{k-1|k-1}^{\xi}\bar{A}_{k-1}^{T} + \bar{Q}_{k-1}$  #Filtering:  
8:  $G_k = P_{k|k-1}^{\xi}\bar{C}_k^T \left(R_k + \bar{C}_k P_{k|k-1}^{\xi}\bar{C}_k^T\right)^{-1}$   
9:  $\hat{\xi}_{k|k} = \hat{\xi}_{k|k-1} + G_k \left(y_k - \bar{C}_k \hat{\xi}_{k|k-1}\right)$   
10:  $P_{k|k}^{\xi} = (I - G_k \bar{C}_k) P_{k|k-1}^{\xi}$   
11: end for #Smoothing:  
12: for  $k = K - 1, K - 2, ..., 1$  do  
13:  $K_k = P_{k|k}^{\xi}\bar{A}_k^T \left(P_{k+1|k}^{\xi}\right)^{-1}$   
14:  $P_{k|K}^{\xi} = P_{k|k}^{\xi} + K_k \left(\hat{\xi}_{k+1|K} - P_{k+1|k}^{\xi}\right) K_k^T$   
15:  $\hat{\xi}_{k|K} = \hat{\xi}_{k|k} + K_k \left(\hat{\xi}_{k+1|K} - \bar{A}_k \hat{\xi}_{k|k}\right)$   
16:  $P_{k+1,k|K}^{\xi} = [P_{k+1|K}^{x} P_{k+1|K}^{x}]^T K_k^T$   
17: Compute  $\hat{u}_{k|K}$  and  $P_{k|K}^{u}$  using (30) from  $\hat{\xi}_{k|K}$  and  $P_{k|K}^{\xi}$   
18: end for #M-step:  
19:  $\gamma_k^{(r)} = \text{Diag} \left\{\hat{u}_{k|K} \hat{u}_{k|K}^T + P_{k|K}^u\right\}$ , for  $k = 1, 2, ..., K$   
20:  $\epsilon = \sum_{k=1}^K \|\gamma_k^{(r)} - \gamma_k^{(r-1)}\|_2^2$   
21:  $r = r + 1$   
22: end while 23: Compute  $\{\hat{x}_{k|K}, \hat{u}_{k|K}\}_{k=1}^K$  using (30) from  $\hat{\xi}_{k|K}$  Outputs:  $\{\hat{x}_{k|K}\}_{k=1}^K$  and  $\{\hat{u}_{k|K}\}_{k=1}^K$  using (30) from  $\hat{\xi}_{k|K}$ 

## **Proposition 2:**

For a given set of inputs, let  $\{\gamma_k^{(r)}: k=1,\ldots,K\}_{r=1}^{\infty}$  be the sequence generated by Algorithm 4 such that  $\|\gamma_k^{(r)}\|_2 < \infty$ . Then, the sequence  $\{\log p(\boldsymbol{Y}_1^K \mid \{\gamma_k\}_{k=1}^K\}_{r=1}^{\infty} \text{ is monotonically nondecreasing and converges to a finite limit in } <math>\mathbb{R}$ .

Proof:

For notational brevity, we denote the observations by  $\boldsymbol{y} = \boldsymbol{Y}_1^K$ , the hidden variables by  $\boldsymbol{z} = \{\boldsymbol{X}_1^K, \boldsymbol{U}_1^K\}$ , and the hyperparameters by  $\boldsymbol{\gamma} = \{\gamma_k\}_1^K$ . In the rth iteration of the SBL algorithm, we derive

$$\begin{split} \boldsymbol{\gamma}^{(r+1)} &= \arg\max_{\boldsymbol{\gamma}} \mathcal{Q}(\boldsymbol{\gamma} \mid \boldsymbol{\gamma}^{(r)}) \\ &= \arg\max_{\boldsymbol{\gamma}} \mathbb{E}_{\boldsymbol{z} \sim p(\boldsymbol{z} \mid \boldsymbol{y}, \boldsymbol{\gamma}^{(r)})} [\log p(\boldsymbol{y}, \boldsymbol{z} \mid \boldsymbol{\gamma})]. \end{split}$$

Then, by the properties of the EM algorithm [41, Section 3.2], the ML objective function in (28) is monotonically increasing,

$$\log p(\boldsymbol{y} \mid \boldsymbol{\gamma}^{(r)}) \le \log p(\boldsymbol{y} \mid \boldsymbol{\gamma}^{(r+1)}).$$

Further, (1) and (2), we derive

$$\tilde{\boldsymbol{y}}_K = \tilde{\boldsymbol{f}}(\boldsymbol{z}, \boldsymbol{W}_1^{K-1}) + \tilde{\boldsymbol{v}}_K,$$

where  $\tilde{\boldsymbol{y}}_K = \begin{bmatrix} \boldsymbol{y}_1^\mathsf{T} & \boldsymbol{y}_2^\mathsf{T} & \dots & \boldsymbol{y}_K^\mathsf{T} \end{bmatrix}^T \in \mathbb{R}^{Kp}$  denotes the concatenated measurement vector. Likewise,  $\tilde{\boldsymbol{v}}_K \in \mathbb{R}^{Kp}$  stacks the measurement noise. It is Gaussian with zero mean and covariance  $\boldsymbol{R} = \operatorname{Blkdiag}(\boldsymbol{R}_1, \dots, \boldsymbol{R}_K)$ , where the operator  $\operatorname{Blkdiag}(\cdot)$  denotes a block diagonal matrix with its arguments as its block diagonal entries. The term  $\tilde{\boldsymbol{f}}(\boldsymbol{z}, \boldsymbol{W}_1^{K-1})$  is a linear function of the initial state, inputs, and process noise. It is also Gaussian distributed with zero mean and covariance denoted by  $\boldsymbol{\Delta}(\boldsymbol{\gamma}) \in \mathbb{R}^{Kp \times Kp}$ , i.e., it is a function of  $\boldsymbol{\gamma}$ . Then, we derive

$$\begin{split} \log p(\boldsymbol{y} \mid \boldsymbol{\gamma}^{(r)}) &= -\frac{Kp}{2} \log(2\pi) - \frac{1}{2} \log |\boldsymbol{\Delta}(\boldsymbol{\gamma}) + \boldsymbol{R}| \\ &- \frac{1}{2} \tilde{\boldsymbol{y}}_K^\mathsf{T} \left(\boldsymbol{\Delta}(\boldsymbol{\gamma}) + \boldsymbol{R}\right)^{-1} \tilde{\boldsymbol{y}}_K \\ &\leq -\frac{Kp}{2} \log(2\pi) - \frac{1}{2} \log |\boldsymbol{R}| \,, \end{split}$$

since  $\Delta(\gamma)$  is positive semi-definite and bounded. Moreover, since R is positive definite, the monotonically nondecreasing sequence  $\{\log p(Y_1^K \mid \{\gamma_k\}_{k=1}^K\}_{r=1}^\infty \text{ is bounded from above and the sequence converges to a finite limit point.}$   $Remark: \text{ When all the inputs are jointly sparse, we can use a common prior } u_k \sim \mathcal{N}(0, \operatorname{Diag}\{\gamma\}), \text{ i.e., } \gamma_k = \gamma \text{ for } k=1,2,\ldots,K. \text{ Then, the SBL-RKS for joint sparse input recovery is identical to Algorithm 4 except for Steps 5 and 19. The noise covariance in Step 5 changes to$ 

$$\bar{\boldsymbol{Q}}_k = \begin{bmatrix} \boldsymbol{Q}_k & \mathbf{0} \\ \mathbf{0} & \operatorname{Diag} \{ \boldsymbol{\gamma} \} \end{bmatrix}, k = 1, 2, \dots, K.$$

Similarly, the M-step in Step 19 is modified as

$$\boldsymbol{\gamma}^{(r+1)} = \frac{1}{K} \sum_{k=1}^{K} \operatorname{Diag} \left\{ \hat{\boldsymbol{u}}_{k|K} \hat{\boldsymbol{u}}_{k|K}^{\mathsf{T}} + \boldsymbol{P}_{k|K}^{\boldsymbol{u}} \right\}.$$

The modified SBL-RKS for jointly sparse inputs is referred to as multiple measurement vector SBL-RKS (MSBL-RKS).

The advantage of SBL is that the hierarchical Gaussian priors together with EM updates automatically learn the penalty weights associated with the relevant nonzero entries of the sparse vector. In contrast, regularization-based approaches require manual tuning of penalty parameters [42].

## B. Variational Bayesian Robust Kalman Smoothing

In the variational Bayesian inference (VBI) approach, we employ a two-stage hierarchical prior. Specifically, we assume  $\beta_k(i) \sim \operatorname{Gamma}(a,b)$ , where the precision hyperparameter  $\beta_k(i) = 1/\gamma_k(i)$  in (27), and  $\operatorname{Gamma}(a,b)$  is the Gamma distribution with shape parameter a>0 and rate parameter b>0, i.e.,

$$p(\beta_k) = \prod_{i=1}^{m} \Gamma^{-1}(a) \ b^a(\beta_k(i))^{a-1} \exp(-b\beta_k(i)).$$
 (31)

The VB-RKS algorithm estimates the set of unknown parameters  $\mathcal{Z} = \left\{ \boldsymbol{X}_1^K, \boldsymbol{U}_1^K, \left\{ \boldsymbol{\beta}_k \right\}_{k=1}^K \right\}$  as the mean of their posterior distribution. However, the posterior distribution computation is intractable, and we approximate it using a family of factorized distributions:

$$p(\mathcal{Z}|\boldsymbol{Y}_1^K) \approx q(\mathcal{Z}) = \prod_{k=1}^K q_k^{\boldsymbol{x}}(\boldsymbol{x}_k) q_k^{\boldsymbol{u}}(\boldsymbol{u}_k) q_k^{\boldsymbol{\beta}}(\boldsymbol{\beta}_k),$$

where  $q_k^{\boldsymbol{x}}(\cdot), q_k^{\boldsymbol{u}}(\cdot)$ , and  $q_k^{\boldsymbol{\beta}}(\cdot)$  are the marginal distributions of the latent variables  $\boldsymbol{x}_k, \boldsymbol{u}_k$ , and  $\boldsymbol{\beta}_k$ , respectively.

We seek the optimal distribution  $q(\mathcal{Z})$  that minimizes the Kullback-Leibler (KL) divergence

$$\mathcal{F}(q(\mathcal{Z})) = \mathrm{KL}(q(\mathcal{Z}) || p(\mathcal{Z} | \boldsymbol{y})).$$

Now, the optimal marginal distribution is given by [43]

$$\ln q_k^{\boldsymbol{x}}(\boldsymbol{x}_k) \propto \mathbb{E}_{q(\boldsymbol{\mathcal{Z}} \backslash \boldsymbol{x}_k)} \left\{ \ln p \left( \boldsymbol{\mathcal{Z}}, \boldsymbol{Y}_1^K \right) \right\} \ln q_k^{\boldsymbol{u}}(\boldsymbol{u}_k) \propto \mathbb{E}_{q(\boldsymbol{\mathcal{Z}} \backslash \boldsymbol{u}_k)} \left\{ \ln p \left( \boldsymbol{\mathcal{Z}}, \boldsymbol{Y}_1^K \right) \right\},$$

where  $\infty$  denotes the equality up to an additive constant and  $\mathbb{E}_{q(\mathcal{Z} \setminus \boldsymbol{x}_k)}\{\cdot\}$  and  $\mathbb{E}_{q(\mathcal{Z} \setminus \boldsymbol{u}_k)}\{\cdot\}$  are the expectations with respect to all the latent variables except  $\boldsymbol{x}_k$  and  $\boldsymbol{u}_k$ , respectively. Further, we recall that

$$p(\mathcal{Z}, \boldsymbol{Y}_1^K) = \prod_{k=1}^K p(\boldsymbol{y}_k | \boldsymbol{x}_k, \boldsymbol{u}_k) p(\boldsymbol{x}_k | \boldsymbol{x}_{k-1}, \boldsymbol{u}_{k-1}) \times p(\boldsymbol{u}_k | \boldsymbol{\beta}_k) p(\boldsymbol{\beta}_k)$$

where  $p(u_k|\beta_k)$  and  $p(\beta_k)$  are given by (27) with  $\beta_k(i) = 1/\gamma_k(i)$ , and (31), respectively. Consequently, we arrive at

$$egin{aligned} & \ln q_k^{oldsymbol{x}}(oldsymbol{x}_k) & \propto \|oldsymbol{y}_k - oldsymbol{C}_k oldsymbol{x}_k - oldsymbol{D}_k \langle oldsymbol{u}_k 
angle \|_{oldsymbol{Q}_k}^2 \ & + \|oldsymbol{x}_k - oldsymbol{A}_{k-1} \langle oldsymbol{x}_{k-1} 
angle - oldsymbol{B}_{k-1} \langle oldsymbol{u}_k 
angle \|_{oldsymbol{Q}_k}^2 \ & + \|oldsymbol{x}_k - oldsymbol{A}_{k-1} \langle oldsymbol{x}_{k-1} 
angle - oldsymbol{B}_{k-1} \langle oldsymbol{u}_{k-1} 
angle \|_{oldsymbol{Q}_{k-1}}^2, \end{aligned}$$

where  $\langle \cdot \rangle$  denotes the mean of a random variable following the marginal distribution  $q(\cdot)$ . Hence, the marginal distribution  $q_k^{\boldsymbol{x}}(\boldsymbol{x}_k)$  is Gaussian. Its mean can be computed by setting the gradient with respect to  $\boldsymbol{x}_k$  to 0, leading to

$$\langle \boldsymbol{x}_{k} \rangle = \boldsymbol{P}_{k}^{\boldsymbol{x}} \left[ \boldsymbol{C}_{k}^{\mathsf{T}} \boldsymbol{R}_{k}^{-1} \boldsymbol{y}_{k} + \boldsymbol{Q}_{k-1}^{-1} \boldsymbol{B}_{k-1} \langle \boldsymbol{u}_{k-1} \rangle - \left( \boldsymbol{C}_{k}^{\mathsf{T}} \boldsymbol{R}_{k}^{-1} \boldsymbol{D}_{k} + \boldsymbol{A}_{k}^{\mathsf{T}} \boldsymbol{Q}_{k}^{-1} \boldsymbol{B}_{k} \right) \langle \boldsymbol{u}_{k} \rangle + \boldsymbol{Q}_{k-1}^{-1} \boldsymbol{A}_{k-1} \langle \boldsymbol{x}_{k-1} \rangle + \boldsymbol{A}_{k}^{\mathsf{T}} \boldsymbol{Q}_{k}^{-1} \langle \boldsymbol{x}_{k+1} \rangle \right], \quad (32)$$



where we define

$$\boldsymbol{P}_{k}^{\boldsymbol{x}} = \left(\boldsymbol{C}_{k}^{\mathsf{T}} \boldsymbol{R}_{k}^{-1} \boldsymbol{C}_{k} + \boldsymbol{Q}_{k-1}^{-1} + \boldsymbol{A}_{k}^{\mathsf{T}} \boldsymbol{Q}_{k}^{-1} \boldsymbol{A}_{k}\right)^{-1}.$$
 (33)

Similarly, the marginal distribution of  $u_k$  is computed as

$$\begin{split} & \ln q_k^{m{u}}(m{u}_k) \propto \mathbb{E}_{q(\mathcal{Z} \setminus m{u}_k)} \left\{ \ln p\left(\mathcal{Z}, m{Y}_1^K \right) \right\} \ & \propto \|m{y}_k - m{C}_k \langle m{x}_k \rangle - m{D}_k m{u}_k\|_{m{R}_k}^2 \ & + \|m{x}_{k+1} - m{A}_k \langle m{x}_k \rangle - m{B}_k m{u}_k\|_{m{Q}_k}^2 + \sum_{i=1}^m m{u}_k(i)^2 \langle m{eta}_k(i) 
angle. \end{split}$$

The mean of the Gaussian distribution  $q_k^{\boldsymbol{u}}(\boldsymbol{x}_k)$  is

$$\langle \boldsymbol{u}_{k} \rangle = \boldsymbol{P}_{k}^{\boldsymbol{u}} \left[ \boldsymbol{D}_{k}^{\mathsf{T}} \boldsymbol{R}_{k}^{-1} \boldsymbol{y}_{k} - \left( \boldsymbol{D}_{k}^{\mathsf{T}} \boldsymbol{R}_{k}^{-1} \boldsymbol{C}_{k} + \boldsymbol{B}_{k}^{\mathsf{T}} \boldsymbol{Q}_{k}^{-1} \boldsymbol{A}_{k} \right) \times \langle \boldsymbol{x}_{k} \rangle + \boldsymbol{B}_{k}^{\mathsf{T}} \boldsymbol{Q}_{k}^{-1} \langle \boldsymbol{x}_{k+1} \rangle \right], \quad (34)$$

where the matrix  $P_k^u$  is

$$\boldsymbol{P}_{k}^{\boldsymbol{u}} = \left(\boldsymbol{D}_{k}^{\mathsf{T}} \boldsymbol{R}_{k}^{-1} \boldsymbol{D}_{k} + \boldsymbol{B}_{k}^{\mathsf{T}} \boldsymbol{Q}_{k}^{-1} \boldsymbol{B}_{k} + \langle \operatorname{diag} \left\{ \boldsymbol{\beta}_{k} \right\} \rangle \right)^{-1}.$$

We use  $x_{K+1} = 0$  for k = K and  $x_0 = 0$  for k = 1 in (32) and (34). Finally,  $q(\pmb{\beta}_k) = \prod_{i=1}^m q(\pmb{\beta}_k(i))$  is a Gamma distribution with mean

$$\langle \boldsymbol{\beta}_k(i) \rangle = \frac{a + 0.5}{b + 0.5 \langle \boldsymbol{u}_k^2(i) \rangle} = \frac{a + 0.5}{b + 0.5 \left[ \langle \boldsymbol{u}_k(i) \rangle^2 + \boldsymbol{P}_k^{\boldsymbol{u}}(i,i) \right]}.$$
(35)

Using (32), (34), and (35), the marginal distribution parameters are iteratively updated until convergence to obtain the approximate posterior distribution. The pseudocode is summarized in Algorithm 5.

## Algorithm 5 Variational Bayesian RKS

Inputs:  $\{y_k, A_k, B_k, C_k, D_k, Q_k, R_k\}_{k=1}^K$ 

**Parameters:**  $r_{\rm max}$  and  $\tilde{r}_{\rm max}$ 

**Initialization:**  $\langle x_k \rangle = 0$ ,  $\langle u_k \rangle = 0$ ,  $\langle \beta_k \rangle = 1$  for k = $1, 2, \ldots, K$ 

- 1: **for**  $r = 1, 2, \dots, r_{\text{max}}$  **do**
- for  $\tilde{r} = 1, 2, \dots, \tilde{r}_{max}$  do
- Compute  $\boldsymbol{x}_k^{(r,\hat{r})} = \langle \boldsymbol{x}_k \rangle$  using (32) for  $k = 1, \dots, K$ Compute  $\boldsymbol{u}_k^{(r,\hat{r})} = \langle \boldsymbol{u}_k \rangle$  using (34) for  $k = 1, \dots, K$ 3:
- 4:
- 5:
- Compute  $\boldsymbol{\beta}_k^{(r)} = \langle \boldsymbol{\beta}_k \rangle$  using (35) for  $k=1,\ldots,K$

7: end for Outputs: 
$$\left\{m{x}_k^{(r, ilde{r})}
ight\}_{k=1}^K$$
 and  $\left\{m{u}_k^{(r, ilde{r})}
ight\}_{k=1}^K$ 

The following result shows that Algorithm 5 monotonically decreases the KL divergence  $\mathcal{F}(\mathcal{Z})$  at each iteration.

#### **Proposition 3:**

For a given set of inputs, the sequence  $\{\mathcal{F}^{(r)}\}_{r=1}^{\infty}$  of KL divergence generated by Algorithm 5 is monotonically nonincreasing and converges to a limit in  $\mathbb{R}$ .

*Proof:* 

From Algorithm 5, let the rth iterate be  $\mathcal{Z}^{(r)}$  $(\boldsymbol{x}^{(r,\tilde{r}_{\max})},\boldsymbol{u}^{(r,\tilde{r}_{\max})},\boldsymbol{\beta}^{(r)})$ . Algorithm 5 proceeds by updating one set of parameters at a time, implying

$$egin{aligned} \mathcal{F}^{(r)} &= \mathcal{F}(\mathcal{Z}^{(r)}) \geq \mathcal{F}(oldsymbol{x}^{(r+1,1)}, oldsymbol{u}^{(r, ilde{r}_{ ext{max}})}, oldsymbol{eta}^{(r)}) \ &\geq \mathcal{F}(oldsymbol{x}^{(r+1,1)}, oldsymbol{u}^{(r+1,1)}, oldsymbol{eta}^{(r)}) \ &\geq \mathcal{F}(oldsymbol{x}^{(r+1, ilde{r}_{ ext{max}})}, oldsymbol{u}^{(r+1, ilde{r}_{ ext{max}})}, oldsymbol{eta}^{(r)}) \ &\geq \mathcal{F}(oldsymbol{x}^{(r+1, ilde{r}_{ ext{max}})}, oldsymbol{u}^{(r+1, ilde{r}_{ ext{max}})}, oldsymbol{eta}^{(r+1)}) = \mathcal{F}^{(r+1)}. \end{aligned}$$

Thus, the sequence  $\{\mathcal{F}^{(r)}\}_{r=1}^{\infty}$  decreases monotonically, and since KL divergence is bounded below by 0, it converges.

We note that our result does not establish that the KL divergence converges to zero, but rather that it converges to a stable value. This indicates that the estimated state and input distributions have stabilized at the closest achievable approximation to the true distribution, though not necessarily identical to it.

Also, when all the inputs are jointly sparse, similar to SBL-RKS, we use a common prior  $u_k \sim \mathcal{N}(\mathbf{0}, \mathrm{Diag}\{\boldsymbol{\beta}\})$ , i.e.,  $\beta_k = \beta$  for k = 1, 2, ..., K. In that case, the VB-RKS for joint sparse input recovery is identical to Algorithm 5 except that (35) in Step 6 changes as follows.

$$\langle \boldsymbol{\beta}(i) \rangle = \frac{a + 0.5}{b + \frac{0.5}{K} \sum_{k=1}^{K} \langle \boldsymbol{u}_k^2(i) \rangle}.$$

We refer to this algorithm as multiple measurement vector VB-RKS (MVB-RKS).

Remark: A special case of our problem is the conventional KF problem when  $u_k = 0 \ \forall k$ . The update steps consist of only two equations (32) and (33) with  $B_k = 0$ ,  $D_k = 0$ . The state update equation is given by,

$$egin{aligned} \langle oldsymbol{x}_k 
angle &= oldsymbol{P}_k^{oldsymbol{x}} igg[ oldsymbol{C}_k^{\mathsf{T}} oldsymbol{R}_k^{-1} oldsymbol{y}_k + oldsymbol{Q}_{k-1}^{-1} oldsymbol{A}_{k-1} igl\langle oldsymbol{x}_{k-1} 
angle \\ &+ oldsymbol{A}_k^{\mathsf{T}} oldsymbol{Q}_k^{-1} \langle oldsymbol{x}_{k+1} 
angle igr]. \end{aligned}$$

VB-RKS allows for simple updates of both states and inputs, whereas the other algorithms typically require running a Kalman filter or an RKS subroutine. Consequently, each iteration of VB-RKS (Step 3, 4, and 6) is computationally faster than those of the other iterative algorithms. However, in simulations (see Section V) we observe that VB-RKS requires significantly more iterations to converge, making it overall slower.

## C. Complexity Comparisons

All four algorithms,  $\ell_1$ -regularized RKS, reweighted  $\ell_2$ regularized RKS, SBL-RKS, and VB-RKS are iterative and uses the RKS algorithm in every iteration. Also, the RKS step is the most computationally complex step in these algorithms and dominates the overall complexity, leading to per-iteration time complexity of all the algorithms as  $\mathcal{O}(K(n^3+m^3+p^3))$  for the versions with and without the joint sparsity assumption. Since the sparsity-driven algorithms consider low-dimensional measurements where  $m \ge$ p, the time complexity further reduces to  $\mathcal{O}(K(n^3+m^3))$ .

Further, we have observed from our simulation results that the other algorithms require a larger number of iterations and a longer time to converge than SBL-RKS (see Table 1). In particular, VB-RKS takes a longer run time owing to the large number of iterations required for convergence. For comparison, we consider the state-of-the-art  $\ell_1$  minimization-based algorithm, referred to as basis pursuit (BP)-RKS (group BP-RKS for the joint support case). BP-RKS is a non-iterative algorithm whose complexity scales as  $\mathcal{O}(K^{\frac{7}{2}}m^{\frac{3}{2}}p^2 + K(n^3 + p^3))$  due to the convex programming optimization using the interior point method. So, our algorithms have low complexity order when the number of iterations is small. The auxiliary space (memory) complexity of all the algorithms is  $\mathcal{O}(p^2 + K(n^2 + m^2))$ . The total time complexities of all the algorithms are summarized in Table 2 for comparison.

#### V. Simulation Results

In this section, we present empirical results to demonstrate the superior performance of the algorithms that exploit sparsity. We choose the state dimension n=30, the input dimension m=100, the output dimension p=20, and the number of time steps K=30. The sparsity level of the input is s = 5, and the locations of s nonzero entries are chosen uniformly at random from the set  $\{1, 2, \ldots, m\}$ . Further, the nonzero entries are drawn independently from a normal distribution  $\mathcal{N}(0, \sigma_u^2)$  with  $\sigma_u = 5$ . For the time-varying support case, we independently choose different supports for each time instant k, and for the jointly sparse case, we use the same support for all values of k. The entries of the time invariant system matrices A, B, C, and D and the initial state  $x_1$  are independently drawn from the standard normal distribution. Also, the process noise covariance Qand the measurement noise covariance R are the identity matrix and  $\sigma_v^2 I$ , respectively. Finally,  $\sigma_v$  is computed from the measurement SNR via the relation SNR =  $s\sigma_u^2/\sigma_v^2$ .

For the above setting, we compare the performance of our algorithms:  $\ell_1$ -regularized RKS, reweighted  $\ell_2$ -regularized RKS, SBL-RKS, and VB-RKS for the time-varying support and jointly sparse cases. We also consider two benchmark approaches: basis pursuit (BP)-RKS and group BP-RKS, which are adapted from the algorithm in [34] (see Appendix B for details), and the RKS algorithm (Algorithm 1). The following metrics are used for comparison: normalized mean squared error (NMSE) in the state and input estimation, false support recovery rate (FSRR) for input estimation, and run time. The FSRR is the sum of the false alarm and missed detection rates of the support estimation. The results in Figs. 1 to 3 and Table 1 compare the algorithms' performance as a function of the measurement dimension p.

## 1) Comparison with RKS and an oracle baseline

Fig. 1 shows different algorithms' NMSE and FSRR performance for the joint sparsity case. We omit the time-varying support case due to space limitations; it can be found in [44, Fig. 1]. From Fig. 1, we infer that the conventional RKS

TABLE 1: Run time comparison of algorithms when n=30, p=20, m=100, K=30, s=5 and SNR is 20 dB

Support	Algorithm	Runtime
Time Varying	RKS	1.2 s
	BP-RKS	73.94 s
	$\ell_1$ -regularized RKS	33.76 s
	reweighted $\ell_2$ -regularized RKS	44.5 s
	SBL-RKS	14.5 s
	VB-RKS	$\sim$ 5 min
	Group BP-RKS	52 s
Jointly	Group $\ell_1$ -regularized RKS	34.6 s
sparse	MSBL-RKS	13.5 s
	MVB-RKS	$\sim$ 5 min

TABLE 2: Complexity comparison of algorithms

Type	Algorithm	Complexity
Non	RKS	$\mathcal{O}(K(n^3 + m^3 + p^3))$
Iterative	BP-RKS	$\mathcal{O}(K^{\frac{7}{2}}m^{\frac{3}{2}}p^2 + K(n^3 + p^3))$
	$\ell_1$ -RKS	$\mathcal{O}(r_{\max}K(n^3+m^3+p^3))$
Iterative	$\ell_2$ -RKS	$\mathcal{O}(r_{\max}K(n^3+m^3+p^3))$
	SBL-RKS	$\mathcal{O}(r_{\max}K(n^3+m^3+p^3))$
	VB-RKS	$\mathcal{O}(r_{\max}\tilde{r}_{\max}K(n^3+m^3+p^3))$

algorithm has poor NMSE performance compared to the sparsity-driven approaches. This underscores the importance of exploiting sparsity to achieve low NMSE. RKS requires p>m for the inverse to exist in Step 2 of Algorithm 1. Naturally, the algorithm fails in the low-dimensional measurement regime. Also, the NMSE of RKS is comparable to that of the sparsity-driven algorithms only when p>m, but the latter outperform RKS even in that regime.

We also compare performance against an oracle RKS baseline that assumes the knowledge of input supports. With the support set  $S_k$  of  $u_k$  known, the system can be represented with reduced matrices  $A_k$ ,  $B_k$ ,  $(C_k)_{S_k}$ , and  $(D_k)_{S_k}$ , driven by non-sparse inputs  $(u_k)_{S_k}$ . Then, oracle RKS estimates  $(u_k)_{S_k}$ . As expected, oracle RKS outperforms all other algorithms in both state and input NMSE performance, as shown in Fig. 1a and Fig. 1b, respectively. However, SBL-RKS and VB-RKS approach oracle-level performance in state estimation in the high-p regime.

#### 2) Comparison of Sparsity-driven Algorithms

Fig. 1 and Table 1 show that the SBL-RKS and VB-RKS algorithms outperform BP-RKS and regularized RKS in terms of NMSE in both states and input estimation, FSRR, and run time, in line with our arguments in Section IV.C. VB-RKS performs similar to SBL-RKS, except for runtime, which is much higher for the former algorithm.

#### 3) Time-varying Support and Joint Sparsity

Fig. 2 and Table 1 show different algorithms' performances for the time-varying support and joint sparsity cases. Fig. 2

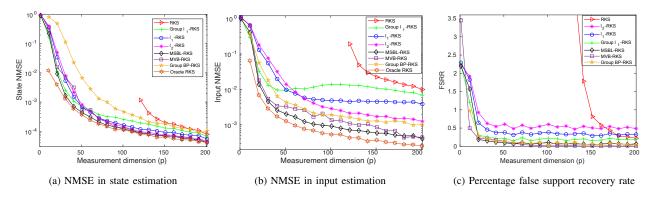


FIGURE 1: Performance comparison of our sparse recovery algorithms and RKS as a function of measurement dimension p when the control inputs are jointly sparse with n = 30, m = 100, K = 30, s = 5, and SNR = 20 dB.

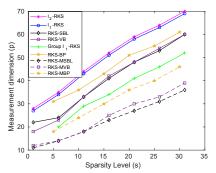
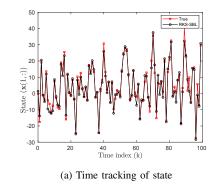


FIGURE 2: Phase transition diagram for our sparse recovery algorithms with  $n=30,\ m=100,\ K=30,$  and SNR = 20 dB.



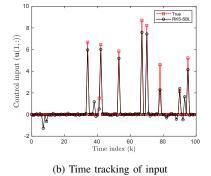


FIGURE 3: Time domain tracking performance of SBL-RKS with  $n=30,\ m=100,\ p=20,\ K=100,\ s=5,$  and SNR =20 dB.

plots the minimum value of the number of measurements prequired for 90% recovery accuracy (i.e., successful recovery of the sparse signals in 90% of the random experiments). Here, a sparse vector is said to be successfully recovered if the normalized mean square error between the original signal and the recovered signal is below 0.05. For the same number of measurements p, the NMSE in input estimation is better for the joint sparsity-aware algorithms: group BP-RKS, group  $\ell_1$ -regularized RKS, MSBL-RKS, and MVB-RKS. This behavior is because of the additional joint sparse structure exploited by these algorithms. Clearly, the Bayesian RKS algorithms require the least number of measurements, followed by group BP-RKS and group  $\ell_1$ regularized RKS. The joint sparsity-aware algorithms are followed by Bayesian-RKS, BP-RKS, and regularized RKS. The regularized RKS algorithms have similar phase transition curves, but SBL-RKS and VB-RKS require fewer measurements than the regularized RKS.

Finally, Table 1 indicates that joint sparsity-aware algorithms have a shorter run time than their counterparts for the time-varying support case because they have fewer parameters to estimate.

## 4) Time domain tracking of states and inputs

The time domain state and input tracking performance of SBL-RKS is shown in Fig. 3 for visual comparison. Here, the estimate of one entry in the state and input vectors across time obtained using SBL-RKS are depicted in Fig. 3a and in Fig. 3b, respectively. It is clear from Fig. 3b that SBL-RKS is able to recover the spikes in the entries of input vectors accurately. Also, Fig. 3a shows that SBL-RKS tracks the state tightly. We observe similar results for all the other sparsity-aware algorithms, and hence, we omit plotting them to avoid clutter.

#### VI. Conclusion

In this paper, we studied the joint estimation of states and sparse inputs as an observer design problem in an LDS. We developed novel algorithms using fictitious sparsity-promoting priors, integrating sparse signal recovery techniques within the Kalman smoothing framework to create sparsity-aware Kalman smoothers. Our methods include a regularization-based MAP estimation approach and a hierarchical Bayesian learning framework built on Gaussian priors. Empirical results show that exploiting sparsity improves estimation and enables recovery with fewer measurements

than conventional methods. Among the proposed methods, SBL-RKS is preferred for its low complexity and superior accuracy. We also extended our approaches for the jointly sparse input case and demonstrated the efficacy of exploiting additional structures with sparsity. We also guarantee convergence of the cost functions of our algorithms to a local minimum or saddle point. However, convergence of the iterates is not ensured due to the combined Kalman smoothing and ADMM, EM, or VBI framework, which needs future study. Investigating the fundamental limits of sparse recovery in LDS is an interesting direction for further research. Future work can also extend our algorithms to handle system identification and estimation of the system matrices.

## Appendix A

# Proof of Theorem 1: The Kalman Filtering and Smoothing Steps

## A. Prediction and Filtering steps

For detailed derivation of step 2-11 see [44]. The filtered estimates of states and inputs are mathematically equivalent to the estimator derived in [20]; hence, we omit it for brevity, although the approach in [44] differs from that in [20].

#### B. Smoothing steps

We develop the smoothing updates for the case of unknown inputs along the lines of the Kalman smoothing in [45]. The smoothed posterior distribution  $p\left(\boldsymbol{x}_{k},\boldsymbol{u}_{k}\mid\boldsymbol{Y}_{1}^{K}\right)$  is Gaussian with mean  $\hat{\boldsymbol{x}}_{k|K},\hat{\boldsymbol{u}}_{k|K}$  and covariance  $\boldsymbol{P}_{k|K}^{\boldsymbol{\xi}}$  as derived below. The joint posterior distribution is

$$p\left(\boldsymbol{\xi}_{k+1}, \boldsymbol{\xi}_{k} \mid \boldsymbol{Y}_{1}^{K}\right)$$

$$= p\left(\boldsymbol{\xi}_{k} \mid \boldsymbol{\xi}_{k+1}, \boldsymbol{Y}_{1}^{k}\right) p\left(\boldsymbol{\xi}_{k+1} \mid \boldsymbol{Y}_{1}^{K}\right)$$

$$= \frac{p\left(\boldsymbol{\xi}_{k+1} \mid \boldsymbol{\xi}_{k}, \boldsymbol{Y}_{1}^{k}\right) p\left(\boldsymbol{\xi}_{k}, \boldsymbol{Y}_{1}^{k}\right)}{p\left(\boldsymbol{\xi}_{k+1}, \boldsymbol{Y}_{1}^{k}\right)} p\left(\boldsymbol{\xi}_{k+1} \mid \boldsymbol{Y}_{1}^{K}\right)$$

$$= \frac{p\left(\boldsymbol{\xi}_{k+1} \mid \boldsymbol{\xi}_{k}\right) p\left(\boldsymbol{\xi}_{k} \mid \boldsymbol{Y}_{1}^{k}\right)}{p\left(\boldsymbol{\xi}_{k+1} \mid \boldsymbol{Y}_{1}^{k}\right)} p\left(\boldsymbol{\xi}_{k+1} \mid \boldsymbol{Y}_{1}^{K}\right),$$

$$= \frac{p\left(\boldsymbol{\xi}_{k+1} \mid \boldsymbol{\xi}_{k}\right) p\left(\boldsymbol{\xi}_{k} \mid \boldsymbol{Y}_{1}^{k}\right)}{p\left(\boldsymbol{\xi}_{k+1} \mid \boldsymbol{Y}_{1}^{K}\right)},$$

$$(36)$$

where (36) follows since  $\{\boldsymbol{y}_{k+1},\ldots,\boldsymbol{y}_K\}$  are linear combinations of  $\boldsymbol{\xi}_{k+1}$  and other random variables  $(\boldsymbol{U}_{k+1}^K,\boldsymbol{W}_{k+1}^K,\boldsymbol{V}_{k+1}^K)$  which are independent of  $\boldsymbol{\xi}_{k+1}$ . Also, (37) follows due to the Markovian nature of the dynamics in (29).

$$p\left(\boldsymbol{\xi}_{k+1}, \boldsymbol{\xi}_{k} \mid \boldsymbol{Y}_{1}^{K}\right)$$

$$= \frac{p\left(\boldsymbol{x}_{k+1} \mid \boldsymbol{x}_{k}, \boldsymbol{u}_{k}\right) p\left(\boldsymbol{u}_{k+1}\right) p\left(\boldsymbol{\xi}_{k} \mid \boldsymbol{Y}_{1}^{k}\right)}{p\left(\boldsymbol{x}_{k+1} \mid \boldsymbol{Y}_{1}^{k}\right) p\left(\boldsymbol{u}_{k+1}\right)} p\left(\boldsymbol{\xi}_{k+1} \mid \boldsymbol{Y}_{1}^{K}\right)$$
(38)

$$= \frac{p\left(\boldsymbol{x}_{k+1} \mid \boldsymbol{x}_{k}, \boldsymbol{u}_{k}\right) p\left(\boldsymbol{\xi}_{k} \mid \boldsymbol{Y}_{1}^{k}\right)}{p\left(\boldsymbol{x}_{k+1} \mid \boldsymbol{Y}_{1}^{k}\right)} p\left(\boldsymbol{\xi}_{k+1} \mid \boldsymbol{Y}_{1}^{K}\right). \tag{39}$$

Here, (38) follows since  $u_{k+1}$  is assumed to be independent of  $x_k$ ,  $x_{k+1}$ ,  $u_k$  and all past observations  $\{y_1, \dots, y_k\}$ .

Taking the logarithm of each term in (39), we have

$$\begin{split} \log p\left(\boldsymbol{x}_{k+1} \mid \boldsymbol{x}_{k}, \boldsymbol{u}_{k}\right) &= -\frac{1}{2} \|\boldsymbol{x}_{k+1} - \boldsymbol{A}_{k} \boldsymbol{x}_{k} - \boldsymbol{B}_{k} \boldsymbol{u}_{k}\|_{\boldsymbol{Q}_{k}}^{2} \\ &= -\frac{1}{2} \|\boldsymbol{T} \boldsymbol{\xi}_{k+1} - \tilde{\boldsymbol{A}}_{k} \boldsymbol{\xi}_{k}\|_{\boldsymbol{Q}_{k}}^{2}, \end{split}$$

where we used  $x_{k+1} = T\xi_{k+1}$  with T and  $A_k$  defined in Algorithm 1. Similarly, we can show that

$$\begin{split} \log p\left(\boldsymbol{x}_{k+1} \mid \boldsymbol{Y}_{1}^{k}\right) &= -\frac{1}{2} \|\boldsymbol{x}_{k+1} - \hat{\boldsymbol{x}}_{k+1|k}\|_{\boldsymbol{P}_{k+1|k}}^{2_{\boldsymbol{x}}} \\ &= -\frac{1}{2} \|\boldsymbol{T}\left(\boldsymbol{\xi}_{k+1} - \hat{\boldsymbol{\xi}}_{k+1|k}\right)\|_{\boldsymbol{P}_{k+1|k}}^{2_{\boldsymbol{x}}} \\ \log p\left(\boldsymbol{\xi}_{k} \mid \boldsymbol{Y}_{1}^{k}\right) &= -\frac{1}{2} \|\boldsymbol{\xi}_{k} - \hat{\boldsymbol{\xi}}_{k|k}\|_{\boldsymbol{P}_{k|k}}^{2_{\boldsymbol{\xi}}} \\ \log p\left(\boldsymbol{\xi}_{k+1} \mid \boldsymbol{Y}_{1}^{K}\right) &= -\frac{1}{2} \|\boldsymbol{\xi}_{k+1} - \hat{\boldsymbol{\xi}}_{k+1|K}\|_{\boldsymbol{P}_{k+1|K}}^{2_{\boldsymbol{\xi}}}. \end{split}$$

Using the above four relations in (39), we derive

$$\log p(\boldsymbol{\xi}_{k+1}, \boldsymbol{\xi}_{k} \mid \boldsymbol{Y}_{1}^{K}) = -\frac{1}{2}\boldsymbol{\xi}_{k+1}^{\mathsf{T}} \Big[ \boldsymbol{T}^{\mathsf{T}} \boldsymbol{Q}_{k}^{-1} \boldsymbol{T} \\ -\boldsymbol{T}^{\mathsf{T}} \left( \boldsymbol{P}_{k+1|k}^{\boldsymbol{x}} \right)^{-1} \boldsymbol{T} + \left( \boldsymbol{P}_{k+1|K}^{\boldsymbol{\xi}} \right)^{-1} \Big] \boldsymbol{\xi}_{k+1} \\ + \frac{1}{2}\boldsymbol{\xi}_{k+1}^{\mathsf{T}} \boldsymbol{T}^{\mathsf{T}} \boldsymbol{Q}_{k}^{-1} \tilde{\boldsymbol{A}}_{k} \boldsymbol{\xi}_{k} + \frac{1}{2}\boldsymbol{\xi}_{k}^{\mathsf{T}} \tilde{\boldsymbol{A}}_{k}^{\mathsf{T}} \boldsymbol{Q}_{k}^{-1} \boldsymbol{T} \boldsymbol{\xi}_{k+1} \\ - \frac{1}{2}\boldsymbol{\xi}_{k}^{\mathsf{T}} \left[ \tilde{\boldsymbol{A}}_{k}^{\mathsf{T}} \boldsymbol{Q}_{k}^{-1} \tilde{\boldsymbol{A}}_{k} + \left( \boldsymbol{P}_{k|k}^{\boldsymbol{\xi}} \right)^{-1} \right] \boldsymbol{\xi}_{k} + \boldsymbol{\xi}_{k}^{\mathsf{T}} \left( \boldsymbol{P}_{k|k}^{\boldsymbol{\xi}} \right)^{-1} \hat{\boldsymbol{\xi}}_{k|k} \\ + \text{linear and constant terms.}$$

Next, we follow the proof technique for deriving the smoothing updates in [45], and we use same notation  $S_{11}$ ,  $S_{12}$ ,  $S_{21}$ ,  $S_{22}$  and  $F_{11}$  as given in the reference which indicate same variables in our context. In our case,  $S_{22}^{-1}$  is

$$egin{aligned} oldsymbol{S}_{22}^{-1} &= \left( ilde{oldsymbol{A}}_{k}^{\mathsf{T}} oldsymbol{Q}_{k}^{-1} ilde{oldsymbol{A}}_{k} + \left( oldsymbol{P}_{k|k}^{oldsymbol{\xi}} 
ight)^{-1} 
ight)^{-1} & = oldsymbol{P}_{k|k}^{oldsymbol{\xi}} - oldsymbol{P}_{k|k}^{oldsymbol{\xi}} ilde{oldsymbol{A}}_{k}^{\mathsf{T}} \left( oldsymbol{Q}_{k} + ilde{oldsymbol{A}}_{k} oldsymbol{P}_{k|k}^{oldsymbol{\xi}} ilde{oldsymbol{A}}_{k}^{oldsymbol{\xi}} - oldsymbol{P}_{k|k}^{oldsymbol{\xi}} ilde{oldsymbol{A}}_{k}^{oldsymbol{\xi}} \left( oldsymbol{P}_{k+1|k}^{oldsymbol{x}} oldsymbol{Q}_{k+1|k}^{oldsymbol{x}} \right)^{-1} ilde{oldsymbol{A}}_{k} oldsymbol{P}_{k|k}^{oldsymbol{\xi}} & = oldsymbol{P}_{k|k}^{oldsymbol{\xi}} - oldsymbol{K}_{k} oldsymbol{P}_{k+1|k}^{oldsymbol{x}} oldsymbol{K}_{k}^{oldsymbol{\xi}}, \end{aligned}$$

where we define

$$oldsymbol{K}_k = oldsymbol{P}_{k|k}^{oldsymbol{\xi}} ilde{oldsymbol{A}}_k^{\mathsf{T}} \left( oldsymbol{P}_{k+1|k}^{oldsymbol{x}} 
ight)^{-1}.$$

Similarly,  $S_{21} = -\tilde{A}_k^\mathsf{T} Q_k^{-1} T$ , and  $S_{22}^{-1} S_{21}$  simplifies to  $S_{22}^{-1} S_{21} = -\left(P_{k|k}^{\boldsymbol{\xi}} - P_{k|k}^{\boldsymbol{\xi}} \tilde{A}_k^\mathsf{T} \left(Q_k + \tilde{A}_k P_{k|k}^{\boldsymbol{\xi}} \tilde{A}_k^\mathsf{T}\right)^{-1} \right)$   $= -P_{k|k}^{\boldsymbol{\xi}} \tilde{A}_k^\mathsf{T} \left(I_n - \left(Q_k + \tilde{A}_k P_{k|k}^{\boldsymbol{\xi}} \tilde{A}_k^\mathsf{T}\right)^{-1} \right)$   $= -P_{k|k}^{\boldsymbol{\xi}} \tilde{A}_k^\mathsf{T} \left(I_n - \left(Q_k + \tilde{A}_k P_{k|k}^{\boldsymbol{\xi}} \tilde{A}_k^\mathsf{T}\right)^{-1} \right)$   $= -P_{k|k}^{\boldsymbol{\xi}} \tilde{A}_k^\mathsf{T} \left(Q_k + \tilde{A}_k P_{k|k}^{\boldsymbol{\xi}} \tilde{A}_k^\mathsf{T}\right)^{-1} T$   $= -P_{k|k}^{\boldsymbol{\xi}} \tilde{A}_k^\mathsf{T} \left(P_{k+1|k}^x\right)^{-1} T = -K_k T.$ 

Hence, the covariance update is given by (see (19) of [45]),

$$\begin{split} & \boldsymbol{P}_{k|K}^{\boldsymbol{\xi}} = \boldsymbol{S}_{22}^{-1} + \boldsymbol{S}_{22}^{-1} \boldsymbol{S}_{21} \boldsymbol{F}_{11}^{-1} \boldsymbol{S}_{12} \boldsymbol{S}_{22}^{-1} \\ & = \left( \boldsymbol{P}_{k|k}^{\boldsymbol{\xi}} - \boldsymbol{K}_{k} \boldsymbol{P}_{k+1|k}^{\boldsymbol{x}} \boldsymbol{K}_{k}^{\mathsf{T}} \right) + \left( -\boldsymbol{K}_{k} \boldsymbol{T} \right) \boldsymbol{P}_{k+1|K}^{\boldsymbol{\xi}} \left( -\boldsymbol{K}_{k} \boldsymbol{T} \right)^{\mathsf{T}} \\ & = \boldsymbol{P}_{k|k}^{\boldsymbol{\xi}} + \boldsymbol{K}_{k} \left( \boldsymbol{P}_{k+1|K}^{\boldsymbol{x}} - \boldsymbol{P}_{k+1|k}^{\boldsymbol{x}} \right) \boldsymbol{K}_{k}^{\mathsf{T}}, \end{split}$$

where we used  $TP_{k+1|K}^{\xi}T^{\mathsf{T}} = P_{k+1|K}^{x}$  in the last step. Hence, we prove Step 14 of Algorithm 1.

Finally, we can compute  $P_{k+1,k|K}^{\xi}$  as (see (20) in [45])

$$\begin{split} \boldsymbol{P}_{k+1,k|K}^{\boldsymbol{\xi}} &= -\boldsymbol{F}_{11}^{-1}\boldsymbol{S}_{12}\boldsymbol{S}_{22}^{-1} = -\boldsymbol{P}_{k+1|K}^{\boldsymbol{\xi}} \left( -\boldsymbol{K}_{k}\boldsymbol{T} \right)^{\mathsf{T}} \\ &= \begin{bmatrix} \boldsymbol{P}_{k+1|K}^{\boldsymbol{x}} \boldsymbol{K}_{k}^{\mathsf{T}} \\ \left( \boldsymbol{P}_{k+1|K}^{\boldsymbol{x}\boldsymbol{u}} \right)^{\mathsf{T}} \boldsymbol{K}_{k}^{\mathsf{T}} \end{bmatrix}. \end{split}$$

We can find the smoothed posterior estimate using [45, Equation (23)], which leads to  $\hat{\boldsymbol{\xi}}_{k|K} = -\boldsymbol{S}_{22}^{-1}\boldsymbol{S}_{21}\hat{\boldsymbol{\xi}}_{k+1|K} + \boldsymbol{S}_{22}^{-1}\left(\boldsymbol{P}_{k|k}^{\boldsymbol{\xi}}\right)^{-1}\hat{\boldsymbol{\xi}}_{k|k}$  and simplifies to

$$\hat{oldsymbol{\xi}}_{k|K} = oldsymbol{K}_k \hat{oldsymbol{x}}_{k+1|K} + \left(oldsymbol{I}_{n+m} - oldsymbol{K}_k ilde{oldsymbol{A}}_k 
ight) \hat{oldsymbol{\xi}}_{k|k}.$$

Thus, we prove Step 15 of Algorithm 1, and it completes proof of all the steps in the algorithm.

## Appendix B

#### **Derivation of BP-RKS and Group BP-RKS**

In this section, we derive two benchmark algorithms inspired by [34], namely BP-RKS for the time-varying support case, and its extension to the joint sparsity case. To derive BP-RKS, we first consider the problem of estimating the initial state  $\boldsymbol{x}_1$  and inputs  $\boldsymbol{U}_1^K$  for the linear system defined in (1) and (2), which can be written as

$$\tilde{\boldsymbol{y}}_{K} = \boldsymbol{O}_{K} \boldsymbol{x}_{1} + \boldsymbol{\Gamma}_{K} \tilde{\boldsymbol{u}}_{K} + \boldsymbol{M}_{K} \tilde{\boldsymbol{w}}_{K-1} + \tilde{\boldsymbol{v}}_{K}, \tag{40}$$

where  $\tilde{\boldsymbol{y}}_K = \begin{bmatrix} \boldsymbol{y}_1^\mathsf{T} & \boldsymbol{y}_2^\mathsf{T} & \dots & \boldsymbol{y}_K^\mathsf{T} \end{bmatrix}^T \in \mathbb{R}^{Kp}$  denotes the concatenated measurement vector. Likewise,  $\tilde{\boldsymbol{u}}_K \in \mathbb{R}^{Km}, \tilde{\boldsymbol{w}}_{K-1} \in \mathbb{R}^{(K-1)n}$ , and  $\tilde{\boldsymbol{v}}_K \in \mathbb{R}^{Kp}$  are obtained by concatenating the inputs, process noise terms and measurement noise terms, respectively. The system matrices,  $\boldsymbol{O}_K \in \mathbb{R}^{Kp \times n}, \; \boldsymbol{\Gamma}_K \in \mathbb{R}^{Kp \times Km} \; \text{and} \; \boldsymbol{M}_K \in \mathbb{R}^{Kp \times (K-1)n}$  in (40) are given by

$$m{O}_{K}\!=\!egin{bmatrix} m{C} \ m{C} \ m{C} \ m{C} \ m{O} \$$

$$\Gamma_{K} = egin{bmatrix} D & 0 & 0 & \cdots & 0 \ CB & D & 0 & \cdots & 0 \ CAB & CB & D & \cdots & 0 \ \vdots & \vdots & \ddots & & \ CA^{K-2}B & CA^{K-3}B & CA^{K-4}B & \cdots & D \end{bmatrix}.$$

For notational simplicity, we present (41) and (42) for the constant system matrices case, i.e.,  $A_k = A$ ,  $B_k = B$ ,  $D_k = C$ , and  $D_k = D$ , k = 1, 2, ..., K. However, the extension to time-varying matrices is straightforward.

Similar to the approach in [34], we first eliminate the initial state term in (40) by multiplying it with  $\Pi$ , the projection matrix onto the orthogonal complement of the column space of the observability matrix  $O_K$ , to get

$$\Pi \tilde{\boldsymbol{y}}_K = \Pi \Gamma_K \tilde{\boldsymbol{u}}_K + \Pi \tilde{\boldsymbol{n}}_K, \tag{43}$$

where the noise term  $\tilde{n}_K = M_K \tilde{w}_{K-1} + \tilde{v}_K$  and the projection matrix is

$$\mathbf{\Pi} = \mathbf{I} - \mathbf{O}_K (\mathbf{O}_K^T \mathbf{O}_K)^{-1} \mathbf{O}_K^T, \tag{44}$$

where  $O_K$  is assumed to have full rank n.

Now, we estimate the states and inputs in two steps, using a Laplacian prior on inputs to encourage sparsity. We first solve for inputs  $\tilde{u}_K$  using (43). Substituting the estimate  $\tilde{u}_K^*$  into (40), we then compute the weighted least square estimate  $x_1^*$  as

$$x_1^* = (O_K^\mathsf{T} Q_{\tilde{n}}^{-1} O_K)^{-1} O_K^\mathsf{T} Q_{\tilde{n}}^{-1} (\tilde{y}_K - \Gamma_K \tilde{u}_K^*).$$
 (45)

Here,  $\tilde{Q}_K$  is the covariance of the noise  $\tilde{n}$  in (40), given by

$$\tilde{\boldsymbol{Q}}_K = \boldsymbol{M}_K \operatorname{Blkdiag}(\boldsymbol{Q}_1, \dots, \boldsymbol{Q}_K) \boldsymbol{M}_K^{\mathsf{T}} + \operatorname{Blkdiag}(\boldsymbol{R}_1, \dots, \boldsymbol{R}_K).$$
 (46)

Having estimated the initial state and inputs, the states  $X_2^K$  can be reconstructed using the Kalman smoothing algorithm.

To estimate the sparse inputs  $\tilde{\boldsymbol{u}}_K$  from (43), we note that the projection step makes  $\Pi\Gamma_K$  rank deficient. This is because from (44), we get  $\mathrm{rank}(\Pi) = Kp - n$  as  $\mathrm{rank}(\boldsymbol{O}_K) = n$ . We further reduce the system in (43) to get linearly independent measurements. Denote the singular value decomposition of the matrix  $\Pi\Gamma_K$  by

$$\mathbf{\Pi}\mathbf{\Gamma}_{K} = \begin{bmatrix} \mathbf{\Psi}_{1} \\ \mathbf{\Psi}_{2} \end{bmatrix}^{\mathsf{H}} \begin{bmatrix} \mathbf{\Lambda} \in \mathbb{R}^{R \times R} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{\Phi}_{1} \\ \mathbf{\Phi}_{2} \end{bmatrix} = \mathbf{\Psi}_{1}^{\mathsf{H}} \mathbf{\Lambda} \mathbf{\Phi}_{1}, \tag{47}$$

where  $\Psi = \begin{bmatrix} \Psi_1^H & \Psi_2^H \end{bmatrix}^H$  is an orthonormal matrix and R is the rank of  $\Pi\Gamma_K$ . When the system matrices have full rank, we have  $R = \min\{Kp - n, Km\}$ . From (43), we derive

$$\boldsymbol{\Psi}\boldsymbol{\Pi}\tilde{\boldsymbol{y}}_{K} = \begin{bmatrix} \boldsymbol{\Psi}_{1}\boldsymbol{\Pi}\tilde{\boldsymbol{y}}_{K} \\ \boldsymbol{\Psi}_{2}\boldsymbol{\Pi}\tilde{\boldsymbol{y}}_{K} \end{bmatrix} = \begin{bmatrix} \boldsymbol{\Lambda}\boldsymbol{\Phi}_{1} \\ \boldsymbol{0} \end{bmatrix} \tilde{\boldsymbol{u}}_{K} + \begin{bmatrix} \boldsymbol{\Psi}_{1}\boldsymbol{\Pi}\tilde{\boldsymbol{n}}_{K} \\ \boldsymbol{\Psi}_{2}\boldsymbol{\Pi}\tilde{\boldsymbol{n}}_{K} \end{bmatrix}. (48)$$

Hence, the reduced system of equations is

$$\Psi_1 \Pi \tilde{\boldsymbol{y}}_K = \Psi_1 \Pi \Gamma_K \tilde{\boldsymbol{u}}_K + \Psi_1 \Pi \tilde{\boldsymbol{n}}_K, \tag{49}$$

where the covariance of the noise component  $\Psi_1\Pi\tilde{n}_K$  is

$$\bar{\boldsymbol{Q}}_K = \boldsymbol{\Psi}_1 \boldsymbol{\Pi} \tilde{\boldsymbol{Q}}_K \boldsymbol{\Pi}^\mathsf{T} \boldsymbol{\Psi}_1^\mathsf{H}. \tag{50}$$

Finally, we multiply the reduced measurements in (49) with the prewhitening matrix  $\bar{Q}_K^{-\frac{1}{2}}$  to make the noise uncorrelated. Hence, the new system of equations is

$$\bar{\boldsymbol{y}}_K = \bar{\boldsymbol{Q}}_K^{-\frac{1}{2}} \boldsymbol{\Psi}_1 \boldsymbol{\Pi} \tilde{\boldsymbol{y}}_K = \bar{\boldsymbol{\Gamma}}_K \tilde{\boldsymbol{u}}_K + \bar{\boldsymbol{n}}_K, \tag{51}$$

## Algorithm 6 Basis Pursuit Robust Kalman Filtering

Inputs:  $\{y_k, A_k, B_k, C_k, D_k, Q_k, R_k\}_{k=1}^K$ 

- 1: Compute  $\tilde{\boldsymbol{y}}_K, \boldsymbol{O}_K, \boldsymbol{M}_K, \boldsymbol{\Gamma}_K$  from (40), (41), (42) and  $\boldsymbol{\Pi}$  using (44)
- 2: Determine  $\Phi_1$  using the singular value decomposition of  $\Pi\Gamma_K$  as given in (47)
- 3: Compute  $\bar{\boldsymbol{Q}}_K$  using (46) and (50)
- 4: Compute  $\bar{\pmb{y}}_K$  and  $\bar{\pmb{\Gamma}}_K$  from (51)
- 5: Set parameter  $\epsilon = \sqrt{R\left(1 + 2\sqrt{\frac{2}{R}}\right)}$
- 6: Solve for inputs  $\dot{\boldsymbol{U}}_{1}^{K}$  using the convex optimization problem (52)
- 7: Calculate the initial state estimate  $x_1$  using (45)
- 8: Set  $x_{1|K} = x_{1|1} = x_1$ Kalman Smoother
- 9: **for** k = 2, ..., K **do** #Prediction:
- 10:  $\hat{\boldsymbol{x}}_{k|k-1} = \boldsymbol{A}_{k-1}\hat{\boldsymbol{x}}_{k-1|k-1} + \boldsymbol{B}_{k}\boldsymbol{u}_{k-1}$
- 11:  $m{P}_{k|k-1}^{m{x}} = m{A}_{k-1} m{P}_{k-1|k-1}^{m{x}} m{A}_{k-1}^{\mathsf{T}} + m{Q}_{k-1}$  #Filtering:

12: 
$$G_k = P_{k|k-1}^x C_k^\mathsf{T} \left( R_k + C_k P_{k|k-1}^x C_k^\mathsf{T} \right)^{-1}$$

- 13:  $\hat{x}_{k|k} = \hat{x}_{k|k-1} + \hat{G}_k \left( y_k D_k u_k C_k \hat{x}_{k|k-1} \right)$
- 14:  $P_{k|k}^{\boldsymbol{x}} = (\boldsymbol{I} \boldsymbol{G}_k \boldsymbol{C}_k) P_{k|k-1}^{\boldsymbol{x}}$
- **15: end for**

#Smoothing:

16: **for** 
$$k = K - 1, K - 2, \dots, 2$$
 **do**

17: 
$$\boldsymbol{K}_{k} = \boldsymbol{P}_{k|k}^{\boldsymbol{x}} \boldsymbol{A}_{k}^{\mathsf{T}} \left( \boldsymbol{P}_{k+1|k}^{\boldsymbol{x}} \right)^{\mathsf{T}}$$

18: 
$$P_{k|K}^x = P_{k|k}^x + K_k \left( P_{k+1|K}^x - P_{k+1|k}^x \right) K_k^\mathsf{T}$$

19: 
$$\hat{oldsymbol{x}}_{k|K} = \hat{oldsymbol{x}}_{k|k} + oldsymbol{K}_k \left( \hat{oldsymbol{x}}_{k+1|K} - oldsymbol{A}_k \hat{oldsymbol{x}}_{k|k} 
ight) - oldsymbol{P}_{k|k}^{oldsymbol{x}} (oldsymbol{I} - oldsymbol{K}_k oldsymbol{A}_k oldsymbol{A}_k^{\mathsf{T}} oldsymbol{Q}_k^{-1} oldsymbol{B}_k oldsymbol{u}_k$$

20: **end for** 

Outputs: 
$$\left\{\hat{\boldsymbol{x}}_{k|K}\right\}_{k=1}^{K}$$
 and  $\left\{\boldsymbol{u}_{k}\right\}_{k=1}^{K}$ 

where matrix  $\bar{\Gamma}_K = \bar{Q}_K^{-\frac{1}{2}} \Psi_1 \Pi \Gamma_K$  and the noise term  $\bar{n}_K = \bar{Q}_K^{-\frac{1}{2}} \Psi_1 \Pi \tilde{n}_K$  follows the standard Gaussian distribution  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ . Under a Laplacian prior on  $\tilde{u}_K$  to promote sparsity, the corresponding MAP estimate is obtained by solving a LASSO problem [46, Section 3.4.3] given by

$$\tilde{\boldsymbol{u}}_K^* = \underset{\tilde{\boldsymbol{u}}_K}{\operatorname{arg\,min}} \|\tilde{\boldsymbol{u}}_K\|_1 \text{ s.t. } \|\bar{\boldsymbol{y}}_K - \bar{\boldsymbol{\Gamma}}_K \tilde{\boldsymbol{u}}_K\|_2 \le \epsilon.$$
 (52)

Here,  $\epsilon>0$  is typically chosen as  $\sqrt{\mathrm{var}(\bar{\boldsymbol{n}}_K)R}\sqrt{1+2\sqrt{\frac{2}{R}}}$  which is slightly larger than  $\sqrt{\mathrm{var}(\bar{\boldsymbol{n}}_K)R}$  (in our case  $\mathrm{var}(\bar{\boldsymbol{n}}_K)=1$ ) [47].

In the second step, we compute the state estimates using the optimal solution  $\tilde{u}_K^*$  obtained by solving (52). Treating

these inputs as the true ones,  $u_k = u_k^*$ , we compute the MAP estimate of the states  $\hat{x}_k$  using the Kalman smoothing algorithm applied to the system (1) and (2). The resulting BP-RKS algorithm is summarized in Algorithm 6.

When control inputs share a common support, we can rearrange  $\tilde{\boldsymbol{u}}_K$  into a block-sparse vector  $\hat{\boldsymbol{u}}_K = [\boldsymbol{u}(1,:)^\mathsf{T} \quad \boldsymbol{u}(2,:)^\mathsf{T} \dots \boldsymbol{u}(m,:)^\mathsf{T}]^\mathsf{T} \in \mathbb{R}^{Km}$  with block length K, where u(i,:) as defined in (26). The corresponding columns of  $\bar{\boldsymbol{\Gamma}}_K$  are also rearranged, which we denote as  $\hat{\boldsymbol{\Gamma}}_K$ , leading to the system of equations  $\bar{\boldsymbol{y}}_K = \hat{\boldsymbol{\Gamma}}_K \hat{\boldsymbol{u}}_K + \bar{\boldsymbol{n}}_K$ . We can now exploit the block sparsity structure by imposing an  $\ell_1/\ell_2$  type penalty [48] on  $\hat{\boldsymbol{u}}_K$  to arrive at the following optimization problem:

$$\hat{\boldsymbol{u}}_{K}^{*} = \underset{\hat{\boldsymbol{u}}_{K}}{\arg\min} \sum_{i=1}^{m} \|\boldsymbol{u}(i,:)\|_{2} \text{ s.t. } \|\bar{\boldsymbol{y}}_{K} - \hat{\boldsymbol{\Gamma}}_{K} \hat{\boldsymbol{u}}_{K}\|^{2} \leq \epsilon.$$
 (53)

The problem (53) can be solved by rewriting it as a second-order cone program using an auxiliary variable t as

$$\min_{\substack{\boldsymbol{t} \in \mathbb{R}^m \\ \hat{\boldsymbol{u}}_K \in \mathbb{R}^{K_m}}} \sum_{i=1}^m \boldsymbol{t}_i \text{ s.t. } \boldsymbol{t}_i \geq \|\boldsymbol{u}(i,:)\|_2 \text{ for } i = 1, \dots, m$$

and 
$$\|\bar{\boldsymbol{y}}_K - \hat{\boldsymbol{\Gamma}}_K \hat{\boldsymbol{u}}_K\|_2^2 \le \epsilon$$
. (54)

The resulting group BP-RKS algorithm is identical to Algorithm 6 except that in Step 6 we solve (54) instead of (52).

### **REFERENCES**

- [1] R. K. Chakraborty, G. Joseph, and C. R. Murthy, "Bayesian learning-based kalman smoothing for linear dynamical systems with unknown sparse inputs," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process.*, 2024, pp. 13431–13435.
- [2] A. Olshevsky, "Minimal controllability problems," *IEEE Trans. Control Network Syst.*, vol. 1, no. 3, pp. 249–258, Sep. 2014.
- [3] M. Siami, A. Olshevsky, and A. Jadbabaie, "Deterministic and randomized actuator scheduling with guaranteed performance bounds," *IEEE Trans. Autom. Control*, vol. 66, no. 4, pp. 1686–1701, Jun. 2020.
- [4] G. Joseph and C. R. Murthy, "Controllability of linear dynamical systems under input sparsity constraints," *IEEE Trans. Autom. Control*, vol. 66, no. 2, pp. 924–931, Apr. 2020.
- [5] G. Joseph, "Controllability of a linear system with nonnegative sparse controls," *IEEE Trans. Autom. Control*, vol. 67, no. 1, pp. 468–473, May 2021.
- [6] N. Wendt, C. Dhal, and S. Roy, "Control of network opinion dynamics by a selfish agent with limited visibility," *IFAC-PapersOnLine*, vol. 52, no. 3, pp. 37–42, 2019.
- [7] G. Joseph, B. Nettasinghe, V. Krishnamurthy, and P. K. Varshney, "Controllability of network opinion in Erdös-Rényi graphs using sparse control inputs," SIAM J. Control Optim., vol. 59, no. 3, pp. 2321–2345, 2021.
- [8] M. Raptis, K. Wnuk, and S. Soatto, "Spike train driven dynamical models for human actions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2010, pp. 2077–2084.
- [9] H. L. Taylor, S. C. Banks, and J. F. McCoy, "Deconvolution with the  $\ell_1$  norm," *Geophys.*, vol. 44, no. 1, pp. 39–52, Jan. 1979.
- [10] M. S. O'Brien, A. N. Sinclair, and S. M. Kramer, "Recovery of a sparse spike time series by  $L_1$  deconvolution," *IEEE Trans. Signal Process.*, vol. 42, no. 12, pp. 3353–3365, 1994.
- [11] A. A. Cárdenas, S. Amin, and S. Sastry, "Research challenges for the security of control systems," *HotSec*, vol. 5, p. 15, Jul. 2008.
- [12] J. Slay and M. Miller, "Lessons learned from the Maroochy water breach," in *Proc. Int. Conf. Crit. Infrastruct. Prot.*, Mar. 2007, pp. 73–82.
- [13] R. Ma, P. Shi, and L. Wu, "Sparse false injection attacks reconstruction via descriptor sliding mode observers," *IEEE Trans. Autom. Control*, vol. 66, no. 11, pp. 5369–5376, 2021.

- [14] H. Yang and R. Ma, "Sparse attack reconstruction for cyber–physical systems via descriptor reduced-order observer," *Nonlinear Anal.: Hy-brid Syst.*, vol. 50, p. 101380, 2023.
- [15] Z. Zhao, Y. Xu, Y. Li, Y. Zhao, B. Wang, and G. Wen, "Sparse actuator attack detection and identification: A data-driven approach," *IEEE Trans. Cybern.*, vol. 53, no. 6, pp. 4054–4064, 2023.
- [16] L. An and G.-H. Yang, "Secure state estimation against sparse sensor attacks with adaptive switching mechanism," *IEEE Trans. Autom. Control*, vol. 63, no. 8, pp. 2596–2603, 2018.
- [17] B. Friedland, "Treatment of bias in recursive filtering," *IEEE Trans. Autom. Control*, vol. 14, no. 4, pp. 359–367, Aug. 1969.
- [18] P. K. Kitanidis, "Unbiased minimum-variance linear state estimation," *Automatica*, vol. 23, no. 6, pp. 775–778, Nov. 1987.
- [19] S. Gillijns and B. De Moor, "Unbiased minimum-variance input and state estimation for linear discrete-time systems," *Automatica*, vol. 43, no. 1, pp. 111–116, Jan. 2007.
- [20] —, "Unbiased minimum-variance input and state estimation for linear discrete-time systems with direct feedthrough," *Automatica*, vol. 43, no. 5, pp. 934–937, May 2007.
- [21] G. Gakis and M. C. Smith, "A limit Kalman filter and smoother for systems with unknown inputs," *Int. J. Control*, vol. 97, no. 3, pp. 532– 542, 2024.
- [22] S. Foucart and H. Rauhut, A Mathematical Introduction to Compressive Sensing. Birkhäuser, 2013.
- [23] G. Joseph and C. R. Murthy, "On the observability of a linear system with a sparse initial state," *IEEE Signal Process. Lett.*, vol. 25, no. 7, pp. 994–998, May 2018.
- [24] —, "Measurement bounds for observability of linear dynamical systems under sparsity constraints," *IEEE Trans. Signal Process.*, vol. 67, no. 8, pp. 1992–2006, 2019.
- [25] W. Dai and S. Yüksel, "Observability of a linear system under sparsity constraints," *IEEE Trans. on Autom. Control*, vol. 58, no. 9, pp. 2372– 2376, Mar. 2013.
- [26] M. B. Wakin, B. M. Sanandaji, and T. L. Vincent, "On the observability of linear systems from random, compressive measurements," in *Proc. IEEE Conf. Decis. Control*, Dec. 2010, pp. 4447–4454.
- [27] B. M. Sanandaji, M. B. Wakin, and T. L. Vincent, "Observability with random observations," *IEEE Trans. Autom. Control*, vol. 59, no. 11, pp. 3002–3007, Aug. 2014.
- [28] N. Vaswani, "Kalman filtered compressed sensing," in *Proc. IEEE Int. Conf. Image Process.*, Oct. 2008, pp. 893–896.
- [29] R. Prasad, C. R. Murthy, and B. D. Rao, "Joint approximately sparse channel estimation and data detection in OFDM systems using sparse Bayesian learning," *IEEE Trans. Signal Process.*, vol. 62, no. 14, pp. 3591–3603, Jun. 2014.
- [30] D. Angelosante, G. B. Giannakis, and E. Grossi, "Compressed sensing of time-varying signals," in *Proc. IEEE Int. Conf. Digit. Signal Process.*, 2009, pp. 1–8.
- [31] A. S. Charles, A. Balavoine, and C. J. Rozell, "Dynamic filtering of time-varying sparse signals via  $\ell_1$  minimization," *IEEE Trans. Signal Process.*, vol. 64, no. 21, pp. 5644–5656, Jul. 2016.
- [32] C. Kurisummoottil Thomas and D. Slock, "Gaussian variational Bayes Kalman filtering for dynamic sparse Bayesian learning," in *Proc. IEEE Intl. Conf. Time Ser. Forecast.*, Mar. 2018.
- [33] M. R. O'Shaughnessy, M. A. Davenport, and C. J. Rozell, "Sparse Bayesian learning with dynamic filtering for inference of time-varying sparse signals," *IEEE Trans. Signal Process.*, vol. 68, pp. 388–403, 2020
- [34] S. Sefati, N. J. Cowan, and R. Vidal, "Linear systems with sparse inputs: Observability and input recovery," in *Proc. IEEE Am. Control Conf.*, 2015, pp. 5251–5257.
- [35] K. Poe, E. Mallada, and R. Vidal, "Necessary and sufficient conditions for simultaneous state and input recovery of linear systems with sparse inputs by  $\ell_1$ -minimization," in *Proc. IEEE Conf. Decis. Control*, 2023, pp. 6499–6506.
- [36] ——, "Invertibility of discrete-time linear systems with sparse inputs," in *Proc. IEEE Conf. Decis. Control*, 2024, pp. 6095–6101.
- [37] D. Wipf and S. Nagarajan, "Iterative reweighted  $\ell_1$  and  $\ell_2$  methods for finding sparse solutions," *IEEE J. Sel. Top. Signal Process.*, vol. 4, no. 2, pp. 317–329, Feb. 2010.
- [38] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends in Machine Learning*, vol. 3, no. 1, pp. 1–122, 2010.

- [39] H. Zou, "The adaptive lasso and its oracle properties," J. Am. Stat. Assoc., vol. 101, no. 476, pp. 1418–1429, 2006.
- [40] D. Angelosante, S. I. Roumeliotis, and G. B. Giannakis, "Lasso-Kalman smoother for tracking sparse signals," in *Proc. Asilomar*, Nov. 2009, pp. 181–185.
- [41] G. McLachlan and T. Krishnan, The EM algorithm and extensions, 2nd ed., ser. Wiley series in probability and statistics. Hoboken, NJ: Wiley, 2008.
- [42] D. Wipf and B. Rao, "Sparse Bayesian learning for basis selection," IEEE Trans. Signal Process., vol. 52, no. 8, pp. 2153–2164, 2004.
- [43] D. Nguyen, "An in-depth introduction to variational Bayes note," SSRN, 2023.
- [44] R. K. Chakraborty, G. Joseph, and C. R. Murthy, "Joint state and input estimation for linear dynamical systems with sparse control," 2023. [Online]. Available: https://arxiv.org/abs/2312.02082
- [45] M. Y. Byron, K. V. Shenoy, and M. Sahani, "Derivation of Kalman filtering and smoothing equations," in *Technical report*. Stanford University, 2004.
- [46] T. Hastie, R. Tibshirani, and J. Friedman, The Elements of Statistical Learning: Data Mining, Inference, and Prediction, 2nd ed. New York, NY: Springer, 2009.
- [47] E. J. Candès and J. Romberg, "\$\ell\_1\$-MAGIC: Recovery of sparse signals via convex programming," 2005, MATLAB routines package.
- [48] Y. C. Eldar and M. Mishali, "Robust recovery of signals from a structured union of subspaces," *IEEE Trans. Inf. Theory*, vol. 55, no. 11, pp. 5302–5316, 2009.



Rupam Kalyan Chakraborty received a B.Tech. in Electronics and Communication Engineering from West Bengal University of Technology in 2016, and M.Tech. degrees in Microwave Engineering from the University of Calcutta and Signal Processing from the Indian Institute of Science, Bangalore in 2018 and 2021, respectively. From 2021–2023, he worked as an Algorithm Design Engineer at Signalchip Innovations on GNSS tracking. Currently, he is a

Ph.D. student in the signal processing systems group at the Delft University of Technology, Netherlands. His research interests include statistical signal processing and compressive sensing.



**Geethu Joseph** received the Ph.D. degree in electrical communication engineering from the Indian Institute of Science, Bangalore, in 2019. She is currently an assistant professor in the signal processing systems group at the Delft University of Technology, Delft, Netherlands. Her research interests include statistical signal processing, network control, and machine learning.



Chandra R. Murthy (S'03–M'06–SM'11–F'23) received Ph.D. degree in Electrical and Computer Engineering from UC San Diego in 2006. He worked at Qualcomm (2000–2002) on WCDMA and 802.11b receivers, and at Beceem Communications (2006–2007) on 802.16e WiMAX. Since 2007, he has been with the Department of Electrical Communication Engineering at IISc Bangalore, where he is now a Professor. His research interests are in the areas of energy harvesting communi-

cations, 5G/6G technologies and compressed sensing. He is a recipient of the MeitY Young Faculty Fellowship from the Govt. of India and the Prof. Satish Dhawan state award for engineering from the Karnataka State Government. He is a fellow of the IEEE.