# The Silence that Speaks: Neural Estimation via Communication Gaps

Shubham Aggarwal                                    sa57@illinois.edu
University of Illinois Urbana Champaign, USA

Dipankar Maity                                      dmaity@charlotte.edu
University of North Carolina at Charlotte, USA

Tamer Başar                                         basar1@illinois.edu
University of Illinois Urbana Champaign, USA

## Abstract

Accurate remote state estimation is a fundamental component of many autonomous and networked dynamical systems, where multiple decision-making agents interact and communicate over shared, bandwidth-constrained channels. These communication constraints introduce an additional layer of complexity, namely, the decision of when to communicate. This results in a fundamental trade-off between estimation accuracy and communication resource usage. Traditional extensions of classical estimation algorithms (e.g., the Kalman filter) treat the absence of communication as 'missing' information. However, silence itself can carry implicit information about the system's state, which, if properly interpreted, can enhance the estimation quality even in the absence of explicit communication. Leveraging this implicit structure, however, poses significant analytical challenges, even in relatively simple systems. In this paper, we propose CALM (Communication-Aware Learning and Monitoring), a novel learning-based framework that jointly addresses the dual challenges of communication scheduling and estimator design. Our approach entails learning not only when to communicate but also how to infer useful information from periods of communication silence. We perform comparative case studies on multiple benchmarks to demonstrate that CALM is able to decode the implicit coordination between the estimator and the scheduler to extract information from the instances of 'silence' and enhance the estimation accuracy.

## 1. Introduction

Remote state estimation plays a foundational role in a wide range of engineering applications, including channel state tracking in cellular networks, spectrum sensing in cognitive radio systems, trajectory prediction in autonomous space exploration, and load forecasting in power grids [1, 2], to name a few. In this work, we focus on the problem of remote estimation of a stochastically evolving discrete-time dynamical system. The overall system (as shown in Fig. 1) comprises a dynamically evolving process along with a team of two cooperating decision-making agents: (1) a collocated scheduler and (2) a remote estimator. The scheduler continuously observes the system state and decides when to communicate this information to the remote estimator, which in turn produces the best possible estimate of the state. The primary objective is to balance the competing goals of maintaining estimation accuracy and minimizing communication overhead.

Departing from classical estimation frameworks, such as the Kalman filter [3] and its numerous variations that treat non-transmissions as missing information, the core contribution of this work is to leverage the implicit information contained in communication gaps (i.e., intervals of silence between transmissions). Our approach is inspired by principles in neuroscience and cognitive systems, where humans can infer intent or internal state even in the absence of sensory input—for example, interpreting pauses or hesitations in speech—making the absence of a signal, a signal in its own right. The key technical challenge in exploiting this insight lies in the fact that the estimation dynamics become tightly coupled with the communication policy, which renders the estimator analytically intractable—even in linear systems for which the Kalman filter was originally developed. Moreover, the scheduler policy is unknown and lacks a closed-form solution.

To address these challenges, we propose CALM, an alternating deep reinforcement learning (DRL) framework based on Proximal Policy Optimization (PPO) that jointly learns both the scheduling policy and a nonlinear estimator. Crucially, the learned estimator is capable of extracting structure from communication silence. Empirical results across several benchmark control domains demonstrate that CALM significantly outperforms traditional linear estimation techniques and heuristic scheduling baselines by effectively exploiting latent information embedded in the no-communication events.

## 2. Related Work

Team decision problems have been extensively studied since the foundational works of Marschak and Radner [4, 5], and serve as a cornerstone for the formulation of distributed decision-making problems involving multiple agents. In such settings, decision-makers (DMs) do not have access to common (centralized) information and collaborate to optimize a common objective. Due to the distributed nature of this setup, coupled with the dynamic nature of the underlying Markov Decision Process (MDP), the information structure becomes critical and significantly influences tractability and optimality [6, 7].

The co-design problem of scheduling and estimation considered in this work can be viewed as a special case of a team decision problem involving two DMs: the scheduler and the estimator. The joint design of control/estimation and scheduling policies has been explored in several prior works [8, 9, 10, 11, 12]. Under a partially nested information structure, it is known that for linear systems, a certainty equivalence property holds, allowing the controller to be designed first (as a function of the conditional estimate of the state), and subsequently the scheduler. However, even in such cases, the conditional estimate may not admit a closed-form expression, making the estimation-scheduling co-design, i.e., the problem in this work, fundamentally hard.

To sidestep analytical intractability, most prior approaches restricted the estimator to be of a linear recursive form [10, 13, 14]. This was often justified by introducing assumptions on the estimator's information set—e.g., partially/completely ignoring the scheduling instants' information [10, 11] or by limiting the problem to scalar systems [9, 10], symmetric policy spaces, or symmetric noise distributions [15, 16]. These assumptions allowed the estimator to be designed in closed form for a given scheduling policy. Attempts to address the informational value of no-communication events have been made in literature [11, 12], where a preliminary complexity analysis was provided. Recent work [17] showed that for a
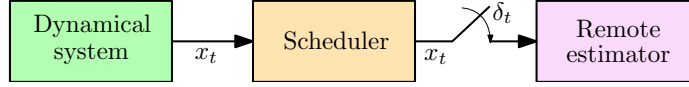
Figure 1: Schematic of a remote estimation system.

stochastic linear system with additive Gaussian noise and one-step delayed communication, silence does not affect estimation, implying that the estimator can remain linear. Nevertheless, the optimal scheduler remains analytically intractable due to the bilinear nature of the estimation error dynamics [18]. The challenge becomes even more pronounced for general nonlinear systems when one attempts to incorporate the full information available at the estimator, including the no-communication instants. Doing so makes the estimator's dynamics nonlinear and non-recursive, breaking the structure that classical solutions rely on, and this has been largely unaddressed in literature.

Contributions: Motivated by the above key technical challenges, in this work we adopt a DRL perspective to study and solve the aforementioned co-design problem. Unlike existing literature, we consider the co-design of the scheduler-estimator policy in remote estimation systems under most general information structures, arbitrary noise distributions and general closed-loop policy classes. To alleviate the consequent intractability yielded within the estimation (and also the scheduling policy) computation, we propose CALM—a DRL-based algorithm for co-designing both the scheduler and the estimator, with the main objective of capturing the latent information present within the periods of silence.

We emphasize that our setup is different from the classical control-oriented neuro-Lyapunov methods, where the objective is that of controller design under full state observability [19, 20, 21] rather than estimator design. Additionally, ours appears to be the first work which deals with general nonlinear stochastic dynamics with arbitrary noise distributions, without restrictions on the noise probability density functions. Finally, our extensive numerical experiments on multiple benchmark control tasks provide insights into how implicit information from silence can be used to infer the underlying structure within the stochastic noise, thereby improving estimation accuracy over existing baselines while keeping communication costs low.

Notations: For an integer $m$, we let $[m] := \{0, 1, 2, \cdots, m\}$. For square matrices $A$ and $B$, we define $\mathrm{diag}(A, B) := \begin{pmatrix} A & 0 \\ 0 & B \end{pmatrix}$. For a positive semi-definite matrix $\Gamma$ and a vector $x$, we define $\|x\|_\Gamma^2 := x^\top \Gamma x$.

## 3. Problem Formulation

Consider a stochastic dynamical system evolving as:

$$x_{t+1} = f(x_t, w_t), \quad t \geq 0, \tag{1}$$

where $x_t \in \mathbb{R}^n$ denotes the state and $w_t \in \mathbb{R}^m$ denotes an i.i.d. noise sample drawn from a zero-mean distribution $\mathcal{P}$, which is assumed to have finite second moments.[1] The initial

---

1. The zero-mean assumption is without any loss of generality. If $\mathbb{E}[w_t] = \bar{w}$, we define $\tilde{w}_t = w_t - \bar{w}$ and $x_{t+1} = f(x_t, w_t) = f(x_t, \bar{w} + \tilde{w}_t) := \tilde{f}(x_t, \tilde{w}_t)$.

state $x_0$ is also sampled from a zero-mean distribution $\mathcal{P}_0$ with finite second moment and is assumed to be independent of the noise distribution $\mathcal{P}$.

The state $x_t$ is actively communicated by a collocated scheduler/sensor to a remote estimator for data logging and analysis, as shown in Fig. 1. However, due to constraints associated with remote communication (e.g., limited power or bandwidth), transmissions must be scheduled in an efficient manner. Naturally, more frequent communication improves estimation accuracy at the expense of higher communication costs, and vice versa. As an illustrative application, consider a Mars rover exploring the Martian surface. It intermittently connects to a NASA base station to transmit its current state. Given strict energy and bandwidth constraints, the question arises: how frequently should the rover communicate its state to ensure sufficient estimation fidelity at the base station while conserving resources?

To pose this question mathematically, let us denote the estimated state at the remote estimator at time $t$ by $\hat{x}_t \in \mathbb{R}^n$, and the scheduling decision by $\delta_t \in \{0,1\}$. Here, $\delta_t = 1$ indicates a communication by the scheduler. Further, let us define the set of (possibly random) scheduling instants up to time $t$ as

$$\mathsf{T}_t := \{t_\ell \mid \ell = 1, \ldots, n_t\}, \tag{2}$$

where $n_t$ is the total number of communications up to time $t$. We assume that the set $\mathsf{T}_t$ is ordered, i.e., $t_1 < t_2 < \cdots < t_{n_t} \leq t$. We have that $\hat{x}_m = x_m$, if and only if, $m \in \mathsf{T}_t$. Henceforth, we will let $\mathsf{T} := \lim_{t \to \infty} \mathsf{T}_t$. Subsequently, we can define the information available to the scheduler and the estimator at time $t$, respectively, as:

$$\begin{aligned}
\mathcal{I}_t^s &:= \{x_k, \delta_k, \hat{x}_k, x_t \mid k \in [t-1]\}, \quad \forall t \geq 1, \\
\mathcal{I}_t^e &:= \{x_m, \delta_k, \hat{x}_k \mid m \in \mathsf{T}_t, k \in [t-1]\}, \ \forall t \geq 1,
\end{aligned} \tag{3}$$

with $\mathcal{I}_0^s := \{x_0, \delta_0\}$ and $\mathcal{I}_0^e := \{\delta_0\}$. Next, we define $\pi$ as the set of measurable scheduling policies according to which decisions $\delta_t$ are made:

$$\pi = \{\mathbf{P}(\delta_t \mid \mathcal{I}_t^s)\}_{\forall t},$$

and the set of measurable estimation policies $\mu$, such that:

$$\mu = \{\mathbf{P}(\hat{x}_t \mid \mathcal{I}_t^e)\}_{\forall t},$$

where $\mathbf{P}(\cdot)$ denotes a Borel-measurable stochastic kernel defined over suitable measurable spaces. The aim is to minimize the following multi-objective cost, expressed as:

$$J(\pi, \mu) := \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t (\|x_t - \hat{x}_t\|_\Gamma^2 + \lambda \delta_t)\right], \tag{4}$$

by co-designing the scheduler-estimator pair $(\delta, \mu)$. In the above, $\Gamma \succeq 0$ denotes a weighting matrix, $\gamma \in (0,1)$ is the discount factor, and $\lambda > 0$ is a scalar that penalizes communication. Thus, we can formally state the problem as follows[2].

---

2. Instead of the uncontrolled system considered in (1), one could indeed consider a controlled dynamical system given by $x_{t+1} = g(x_t, u_t, w_t)$, and subsequently substitute for a stabilizing feedback control $u_t = h(x_t)$ (under suitable conditions of well-posedness to ensure its existence) to arrive at the closed-loop uncontrolled system $x_{t+1} = g(x_t, h(x_t), w_t) =: f(x_t, w_t)$.

**Problem 1** We wish to solve the following:

$$\inf_{\pi,\mu} J(\pi,\mu) \quad \text{subject to} \quad (1),(3). \tag{5}$$

## 4. Scheduler-Estimator Characterization

We begin by stating the following standard result that characterizes the structure of the optimal estimator under a fixed scheduling policy.

**Proposition 1** [11] Suppose that a scheduling policy $\pi$ is fixed. Then, the optimal estimator for Problem 1 is given by the conditional expectation of the state, conditioned on the estimator's information set:

$$\hat{x}_t = \mu^*(\mathcal{I}_t^e) = \mathbb{E}[x_t \mid \mathcal{I}_t^e], \quad \forall t \geq 0. \tag{6}$$

Next, to motivate the technical challenges in the design of (6), and the optimal scheduling policy $\pi^*$, let us consider a simpler exposition of linear systems in the next subsection.

### 4.1. Key Technical Challenges: A Linear System Case Study

Let us begin by considering a linear-time-invariant system, which is a special case of the nonlinear system (1), given as:

$$x_{t+1} = Ax_t + w_t, \quad \forall t \geq 0, \tag{7}$$

where $A \in \mathbb{R}^{n \times n}$ denotes the system matrix. Subsequently, by employing the optimal estimator presented in Proposition 1 and taking the conditional expectation in (7), we obtain the following estimator dynamics:

$$\hat{x}_{t+1} = \begin{cases} A\mathbb{E}[x_t \mid \mathcal{I}_{t+1}^e] + \mathbb{E}[w_t \mid \mathcal{I}_{t+1}^e], & \text{if } \delta_{t+1} = 0, \\ x_{t+1}, & \text{if } \delta_{t+1} = 1, \end{cases} \tag{8}$$

for all $t \geq 0$, with $\hat{x}_0 = (1 - \delta_0)\mathbb{E}[x_0] + \delta_0 x_0$. Note that, when $\delta_{t+1} = 0$, from (3), we have $\mathcal{I}_{t+1}^e = \mathcal{I}_t^e \cup \{\delta_{t+1} = 0\}$.

Define the estimation error as $e_t := x_t - \hat{x}_t$. Then, using (7) and (8), the evolution of the estimation error is given by:

$$e_{t+1} = x_{t+1} - \hat{x}_{t+1} = \begin{cases} Ae_t + \hat{w}_t =: e_{t+1}^o, & \text{if } \delta_{t+1} = 0, \\ 0, & \text{if } \delta_{t+1} = 1, \end{cases} \tag{9}$$

with $e_0 = x_0 - \hat{x}_0$, $\hat{w}_t := w_t - \mathbb{E}[w_t \mid \mathcal{I}_{t+1}^e] + A(\hat{x}_t - \mathbb{E}[x_t \mid \mathcal{I}_{t+1}])$, $e_0^o = x_0 - \mathbb{E}[x_0]$, and we define $e_t^o$ to be the one-step lookahead error at time $t$, assuming no scheduling (i.e., $\delta_t = 0$). Accordingly, the cost function in (4) can be equivalently expressed as:

$$J(\pi, \mu^*) := \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t ((1 - \delta_t)\|e_t^o\|_{\Gamma}^2 + \lambda\delta_t)\right], \tag{10}$$

We can now restate our objective toward solving Problem 1 (for this linear system) as minimizing (10) over $\pi$ subject to (8), (9) and (3). There are, however, two technical challenges associated with the design of a solution to this problem, namely, the computation of the optimal estimator, and the subsequent computation of the scheduling policy.

4.1.1. Technical Challenge #1:

The presence of the conditional expectation terms in (8) renders the estimation dynamics nonlinear, and hence, not expressible in closed form. The challenge becomes even more pronounced for the case of general nonlinear systems. In some cases, however, as presented in the following remarks, this issue can be circumvented but only by imposing additional restrictions over the problem structure, even for the simple case of linear systems.

Remark 2 (Restrictions over the information sets) Existing works (see the dissertation by Molin [10] and the references therein and [18]) simplify the estimator design by not including the scheduling signal $\delta_t$ in the estimator's information set $\mathcal{I}_t^e$. Thus, one disregards the no-communication events (denoted by $\delta_t = 0$), and uses the approximation:

$$\mathbb{E}[x_t \mid \mathcal{I}_{t+1}^e] = \mathbb{E}[x_t \mid \mathcal{I}_t^e \cup \{\delta_{t+1} = 0\}] \overset{(\dagger)}{\approx} \mathbb{E}[x_t \mid \mathcal{I}_t^e] = \hat{x}_t,$$

where ($\dagger$) is precisely where the no-communication events are dropped from the conditional expectation. Additionally, the assumption that $w_t$ is independent across time, and has zero mean, yields $\mathbb{E}[w_t \mid \mathcal{I}_{t+1}^e] = 0$. Consequently, the estimator dynamics in (8) reduce to a piecewise linear form:

$$\hat{x}_{t+1} = \delta_{t+1} x_{t+1} + (1 - \delta_{t+1}) A \hat{x}_t,$$

which can now be computed separately, and recursively for a given scheduling policy.

Remark 3 (Restrictions on noise distribution and policy space) Another common restriction in the literature pertains to that on the noise distribution and the scheduler policy space [15, 16]. Let us denote the domain of the state space where the scheduler decides not to schedule by $\mathcal{D}$. If $\mathcal{D}$ is symmetric about the origin—e.g., $\mathcal{D} := \{z \in \mathbb{R}^n \mid \|z\| \leq a \text{ for some } a > 0\}$—and the noise distribution is also symmetric over $\mathcal{D}$, the estimator in (8) can again be computed using recursive form as presented in Remark 2 [12].

Remark 4 (Restriction to scalar linear systems) We finally remark that linear piecewise estimation (as in Remark 2) is known to be optimal for scalar systems under zero-mean Gaussian noise [8, 9]. In such cases, the absence of communication does not provide additional information for the estimator, and the scheduling policy turns out to be of threshold type, i.e. $\delta = 1$ whenever $|e| \geq \tau$, for a positive scalar $\tau$. However, for multivariate systems, no general result exists, except in certain special cases such as the 'innovation signal'-based scheduling policy [12] or the one-step-delay information transmission pattern [17].

The above discussion (even though involving linear systems) motivates the first technical challenge in estimation design. In this work, we consider the case of general nonlinear systems (and this appears to be the first such study addressing joint scheduler-estimator design in this context) and do not impose any of the simplifying restrictions noted above. In particular, we retain full information in the estimator, allow for general noise distributions, and do not restrict the policy class, for which we will later resort to a learning-based design after presenting the second technical challenge below.

4.1.2. Technical Challenge #2:

Now, we elaborate on the second challenge, which is that of optimal scheduling policy computation. To this end, let us define an MDP $\mathtt{M} := (\mathtt{S}, \mathtt{A}, \mathtt{P}, \mathtt{C})$, where $\mathtt{S}$ is the state space, $\mathtt{A}$ is the action space, $\mathtt{P}$ describes the transition dynamics, and $\mathtt{C}$ defines the per-step cost. In our setting, the MDP state is defined as the lookahead estimation error $e_t^o$, with its dynamics as in (9), and hence $\mathtt{S} = \mathbb{R}^n$. The action space $\mathtt{A} = \{0, 1\}$ corresponds to the binary scheduling decisions with the running cost:

$$\mathtt{C}(e, \delta) = (1 - \delta)\|e^o\|_\Gamma^2 + \lambda\delta.$$

Next, let us fix a scheduling policy $\pi$, and define the state-action value function $Q^\pi(e^o, \delta)$ : $\mathtt{S} \times \mathtt{A} \to \mathbb{R}$ of the MDP:

$$Q^\pi(e^o, \delta) = \mathbb{E}\left[\sum_{t=0}^\infty \gamma^t \mathtt{C}(e_t^o, \delta_t) \mid \pi, e_0^o = e, \delta_0 = \delta\right].$$

Using the Bellman equation, for any time $t$, we write the optimal state-action value function as

$$\begin{aligned} Q(e_t^o, \delta_t) &= \inf_\pi Q^\pi(e_t^o, \delta_t) \\ &= \mathbb{E}[\mathtt{C}(e_t^o, \delta_t) + \gamma \min_{\delta' \in \mathtt{A}} Q(e_{t+1}^o, \delta')]. \end{aligned} \tag{11}$$

Subsequently, one may compute the optimal scheduling policy as

$$\pi^*(\mathcal{I}_t^s) = \arg\min_{\delta_t \in \mathtt{A}} Q(e_t^o, \delta_t). \tag{12}$$

The challenge, however, is that in the given definition, the closed-form expression of the state-action value function remains unknown, making it impossible to compute the scheduling policy exactly.

Remark 5 We remark that there have been multiple attempts to compute the optimal scheduling policy, for the simpler case of linear systems [22, 12, 17, 16]. For instance, the work [17] computes an optimal scheduling policy for a (one-step delayed) linear system only to discover that the policy depends on the value function of the system, and hence, cannot be implemented without appropriate approximations. To alleviate this issue, a recent work [18] attempts to compute the scheduling policy using a deep Q-learning algorithm, albeit by approximating the estimator to its piecewise linear form as discussed earlier.

Building on the preceding discussions on the key challenges, this work is the first to consider the co-design problem for nonlinear systems to explicitly show that communication gaps (i.e., no-scheduling events) can convey implicit information that improves estimation performance in multivariate systems, and without imposing any restrictions as presented before. However, due to the underlying intractability within the estimator, and the coupled scheduling policy design, we resort to a DRL-based co-design framework within an actor-critic architecture: the actor (scheduler deep neural network (DNN)) selects a communication policy, while the critic (estimator DNN) evaluates and updates the state estimate accordingly. We detail this approach in the following section.

## 5. Method & Algorithm

We now describe the alternating training algorithm used to learn both the communication scheduling policy and the state estimator dynamics. The scheduler is modeled using a neural network function approximator which is then trained via PPO with the clipped surrogate objective [23]. The estimator is also modeled using a neural network function approximator, which is trained to predict the system state under sporadic observations. By alternating between training the scheduler (while holding the estimator fixed) and training the estimator (with the scheduler fixed), we enable coordinated learning to balance estimation accuracy and communication cost.

### 5.0.1. PPO subroutine.

We briefly outline the PPO algorithm [23], which serves as a subroutine in our framework. PPO is a first-order reinforcement learning method that addresses instability in traditional policy gradient approaches by limiting large updates via a clipped surrogate objective, thereby offering a simpler and more efficient alternative to TRPO [24].

The foundations of PPO lie within the standard technique of estimating the gradient of expected returns in policy gradient formulation:

$$\hat{g} = \mathbb{E}_t \left[ \nabla_\theta \log \pi_\theta(\delta_t \mid e_t^o) \hat{A}_t^v \right],$$

which corresponds to maximizing the objective function:

$$L^{\mathrm{PG}}(\theta) = \mathbb{E}_t \left[ \log \pi_\theta(\delta_t \mid e_t^o) \hat{A}_t^v \right].$$

Here, $\pi_\theta$ is the current policy parameterized by $\theta$, $\hat{A}_t^v$ is an estimator of the advantage function [23], and $\mathbb{E}_t$ denotes the empirical average taken over a finite number of samples. To ensure stability, PPO introduces a clipped objective:

$$L^{\mathrm{CLIP}}(\theta) = \mathbb{E}_t \left[ \min \left( z_t(\theta) \hat{A}_t^v, \ \mathrm{clip}(z_t(\theta), 1 \pm \epsilon) \hat{A}_t^v \right) \right],$$

where $z_t(\theta) = \frac{\pi_\theta(\delta_t | e_t^o)}{\pi_{\theta_{\mathrm{old}}}(\delta_t | e_t^o)}$ is the policy ratio between the current and the old policy, and $1 \gg \epsilon > 0$ is a hyperparameter. This clipping limits the incentive for overly large updates to remain within a 'trust' region, promoting more stable learning. The resulting algorithm alternates between sampling data using the current policy and optimizing the clipped objective over multiple epochs of stochastic gradient ascent.

### 5.1. Proposed Algorithm: CALM

In light of the above discussion on PPO, we propose CALM: Communication-Aware Learning and Monitoring (presented in Algorithm 1) to jointly construct the scheduler and the estimator. The algorithm alternates between updating the communication scheduling policy via PPO, and subsequently refining the estimator using stochastic gradient descent (SGD). The procedure begins by initializing three neural networks: the estimator network, the policy network, and the value function network. An outer loop (line 1) is then executed, within

which the scheduler is trained (lines 2–15) followed by the estimator training (lines 16–26). To enable the estimator to more effectively capture temporal dynamics, we incorporate an age-of-information (AoI) feature as an additional input during training. Specifically, the AoI is defined as the time elapsed since the last transmission from the scheduler to the estimator. If the current time is $t$, and the most recent communication occurred at time $t_{n_t}$, then AoI $= t - t_{n_t}$, where $n_t$ is defined as in (2). The training is performed in a centralized manner, i.e., each DM maintains a local copy of the trained network of the other DM. During execution, however, both DMs act solely based on the information available to them, as defined in 3. Further, to ensure faster learning, we choose

$$
\mu^*(\mathcal{I}_{t+1}^e) = \begin{cases} f(\hat{x}_t, 0) + \xi_\psi(\hat{x}_t, t - t_{n_t}), & \delta_{t+1} = 0, \\ x_{t+1}, & \delta_{t+1} = 1, \end{cases}
$$

i.e., the estimator DNN (denoted by $\xi_\psi(\cdot)$) only estimates the residual sum: $\mathbb{E}[f(x_t, w_t) \mid \mathcal{I}_{t+1}^e] - f(\hat{x}_t, 0)$.

## 6. Experiments and Analysis

### 6.0.1. Benchmark Details:

We evaluate the proposed CALM framework across three standard control tasks: inverted pendulum stabilization, Van der Pol oscillator, trajectory tracking problem, and a Boeing flight control system. For all experiments, we use the p-mode Gaussian mixture to model the noise distribution, which naturally occurs in diverse mechanical systems such as flight control, human-robot synergies, etc. [25, 26, 27]. Additional implementation and environment-specific details, and signal trajectories corresponding to all experiments are provided in the Appendix.

### 6.1. Results

### 6.1.1. Inverted Pendulum.

For our first experiment, we consider an inverted pendulum system subject to noise drawn from a two-mode Gaussian mixture model (GMM), as shown in (the left subfigure of) Fig. 2. The GMM has component means at $(-3, -3)$ and $(3, 3)$, equal covariance matrices $\text{diag}(0.5, 0.5)$, and respective mixture weights of 0.3 and 0.7. For this example, we also perform an ablation study, where we evaluate performance under a three-mode GMM with component means $(-3, -3)$, $(-5, 4)$, and $(4, 4)$, and covariance matrices $\text{diag}(0.5, 0.5)$, $\begin{pmatrix} 1.0 & 0.8 \\ 0.8 & 1.0 \end{pmatrix}$, and $\begin{pmatrix} 0.6 & -0.3 \\ -0.3 & 0.5 \end{pmatrix}$, with associated mixture weights $(0.6, 0.3, 0.1)$. The communication cost is set to $\lambda = 45$, with further experimental details provided in the Appendix. In Fig. 2, we visualize the scheduling policy learned by CALM as a function of the lookahead estimation error. The resulting decision landscape clearly partitions the state space into two regions: one favoring communication (in red), and the other favoring silence (in cyan).

Silence that speaks: This behavior reflects a form of implicit coordination between the scheduler and the estimator. When the error in Fig. 2 (on the left) is likely drawn

---

**Algorithm 1:** CALM: PPO-driven Alternating Scheduler-Estimator Training Algorithm

---

**Input:** Initial estimator $\xi_{\psi_0}$, scheduling policy $\pi_{\theta_0}$, value function $V_{\phi_0}$, noise distributions $\mathcal{P}, \mathcal{P}_0$, cost weight $\lambda$, rollout length $T$

**Output:** Trained policy $\pi_\theta$ and estimator $\xi_\psi$

for each outer iteration $i = 1$ to $N$ do

    // Train Scheduler using PPO with fixed Estimator

    for each PPO epoch do

        for each trajectory rollout do

            Initialize $x_0 \sim \text{Uniform[-1,1]}$, $\hat{x}_0 \leftarrow \mathbb{E}[x_0]$, $e_0 = x_0 - \hat{x}_0$, $t_{n_t} \leftarrow 0$;

            for $t = 0$ to $T - 1$ do

                Observe error $e_t$;

                Sample action $\delta_t \sim \pi_\theta(e_t)$;

                if $\delta_t = 1$ then

                    $\hat{x}_t \leftarrow x_t$;

                    $t_{n_t} \leftarrow t$;

                end

                Compute reward $r_t = -\left(\|x_t - \hat{x}_t\|_\Gamma^2 + \lambda \cdot \mathbb{I}_{\{\delta_t=1\}}\right)$;

                $x_{t+1} \leftarrow f(x_t, w_t)$;

                $\hat{x}_{t+1} \leftarrow f(\hat{x}_t, 0) + \xi_\psi(\hat{x}_t, t - t_{n_t})$;

            end

            Store $(e_t, \delta_t, \log \pi_\theta(\delta_t|e_t), z_t, V_\phi(e_t))$ and compute advantage estimate $\hat{A}_t^v$;

        end

        Perform PPO update on $\theta$ and $\phi$ using collected trajectories

    end

    // Train Estimator using Fixed Policy

    for each estimator training epoch do

        Initialize $x_0 \sim \text{Unif[-1,1]}$, $\hat{x}_0 \leftarrow \mathbb{E}[x_0]$, $t_{n_t} \leftarrow 0$;

        for $t = 0$ to $T - 1$ do

            Observe error $e_t = x_t - \hat{x}_t$;

            Sample action $\delta_t \sim \pi_\theta(e_t)$;

            if $\delta_t = 1$ then

                $\hat{x}_t \leftarrow x_t$, $t_{n_t} \leftarrow t$;

            end

            Accumulate loss $L$: $L = \gamma L + \|x_t - \hat{x}_t\|_\Gamma^2 + \lambda \cdot \mathbb{I}_{\{\delta_t=1\}}$

            $x_{t+1} \leftarrow f(x_t, w_t)$;

            $\hat{x}_{t+1} \leftarrow f(\hat{x}_t, 0) + \xi_\psi(\hat{x}_t, t - t_{n_t})$;

        end

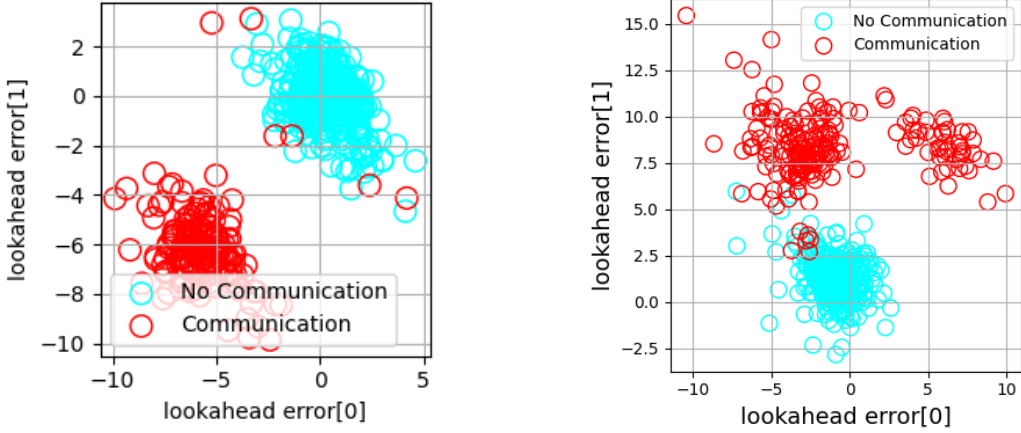        Update estimator parameters $\psi$ using $\nabla_\psi L$;

    end

end

---

Figure 2: Scheduling landscapes for a 2-mode GMM (left) and 3-mode GMM (right) for the inverted pendulum system (red denotes transmissions while cyan denotes silence).

from the positively shifted GMM mode, the scheduler chooses not to transmit (cyan), and the estimator correctly infers this and incorporates the corresponding mode mean into its estimate update. Conversely, when the samples are likely drawn from the negatively shifted mode, the scheduler actively communicates (red), and the estimator uses the actual state. This is a clear instance of "communication via silence" (or a signaling type behavior), where the absence of a transmission conveys informative structure about the underlying noise distribution. Similarly, in the right subfigure in Fig. 2, when samples are drawn from the component in the third quadrant, no communication takes place while it does when sampled from the other two regions.

In contrast, a baseline linear estimator (i.e., the best Kalman filter)—ignorant of this implicit information—simply adds the overall GMM mean without any knowledge of the noise structure, leading to higher overall incurred cost. Specifically, corresponding to Fig. 2 (left), CALM resulted in 168 communications (for $T = 500$) with a total cost of 8355.16 as per (10), while the best linear baseline registered 256 transmissions and incurred a cost of 13386.55. Similarly, for Fig. 2 (right), CALM resulted in 199 communications with a total cost of 10385.45, while the linear baseline registered 174 transmissions and incurred a cost of 11174.44. Results corresponding to the baseline are presented in the Appendix. These results underscore the central premise of our work: silence can carry structured information, which, if leveraged properly, can significantly improve estimation performance.

### 6.1.2. Van der Pol (VdP) oscillator.

Our next case study involves the VdP oscillator, a nonlinear dynamical system of dimension 2. As in the inverted pendulum case, we consider a two-mode GMM with component means at $(-5, -4)$ and $(4, 5)$, both having identical covariance matrices diag$(0.5, 0.5)$. In Fig. 3, we visualize the scheduling landscape generated by CALM as a function of the lookahead estimation error.

Similar to the inverted pendulum case, we again observe the emergence of an agreement between the scheduler and the estimator, resulting in a clear partitioning of the decision
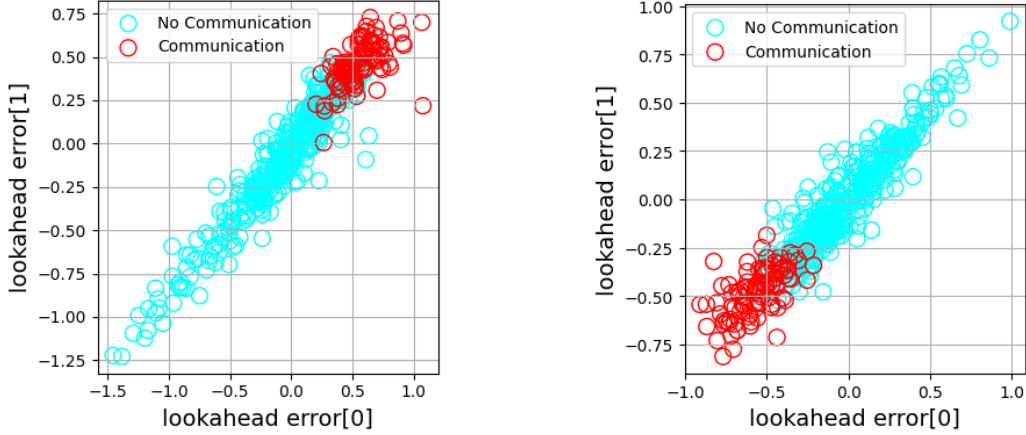
11

Figure 3: Scheduling landscapes for different weights of the 2-modes of GMM for VdP system: weight vector $= (0.3, 0.7)$ on the left and weight vector $= (0.7, 0.3)$ on the right.

space. Further, the positively shifted GMM mode is selected for non-communication in the left subfigure of Fig. 3 due to its higher likelihood (weight of 0.7), thereby minimizing the need for explicit communication while preserving estimation fidelity. As an ablation study, when the GMM weights are reversed, the scheduler-estimator coordination also inverts. This results in a flipped decision landscape, where the previously silent mode now triggers communication and vice versa, as illustrated in the right subfigure of Fig. 3.

### 6.1.3. Robot trajectory tracking.

Our next experiment involves a trajectory tracking system in 2 dimensions, similar to the aforementioned Rover example. For this study, we fix a 4-mode GMM with mean vectors $(-3, -3), (-5, 4), (2, -2), (4, 4)$ and covariance matrices $\mathrm{diag}(0.5, 0.5)$, $\begin{pmatrix} 1 & 0.8 \\ 0.8 & 1 \end{pmatrix}$, $\begin{pmatrix} 0.6 & -0.3 \\ -0.3 & 0.5 \end{pmatrix}$, and $\begin{pmatrix} 0.3 & 0.1 \\ 0.1 & 0.4 \end{pmatrix}$. The weight vector was set to $(0.4, 0.3, 0.2, 0.1)$. In Fig. 4, we investigate how the learned scheduling landscape evolves as a function of the communication cost parameter $\lambda$. For the smallest cost setting, $\lambda = 15$ (top-left), the system has sufficient communication resources to enable communication (indicated by red) for all the four modes. As the cost increases to $\lambda = 30$ (top-right), we observe a no-communication region (indicated in cyan) in the fourth quadrant while communicating regions in the other three quadrants.

At a more restrictive communication cost of $\lambda = 40$ (bottom-left), the scheduler now begins to communicate in only two of the four GMM modes, with increased selectivity under tighter communication constraints. Finally, when the communication cost becomes very high ($\lambda = 70$ (bottom-right)), the budget becomes too limited, and hence, we see a denser cyan region of no-communications compared to the lighter red regions (indicating communication). Further, the communication region also swaps to the first quadrant (from the third quadrant for the case with $\lambda = 40$) since it has a lesser likelihood of being sampled (only 0.1 compared to 0.4 in case of the mode in the third quadrant). This adaptive trade-off highlights the algorithm's ability to allocate communication resources intelligently based on
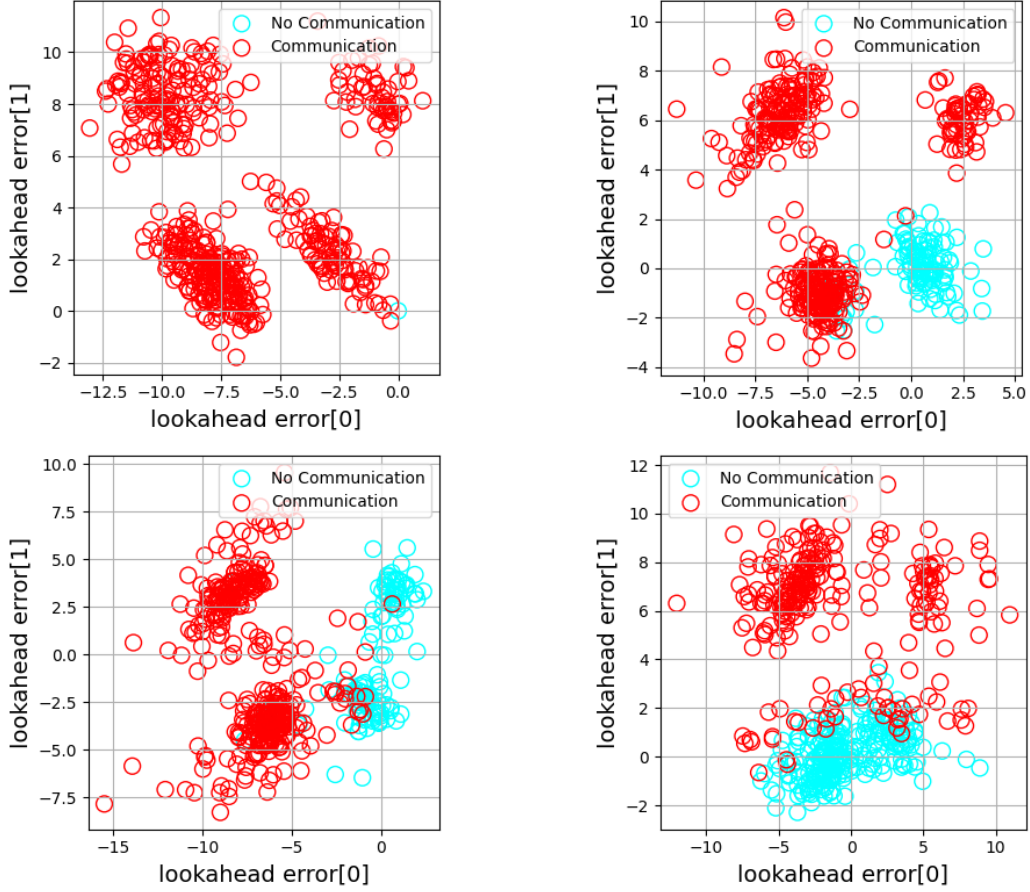
12

Figure 4: Variation of scheduling landscape with communication cost $\lambda$ for the trajectory tracking experiment: $\lambda = 15$ (top left), $\lambda = 30$ (top right), $\lambda = 40$ (bottom-left) and $\lambda = 70$ (bottom right).

both statistical structure (e.g., covariance and probability mass) and cost parameter. This trend can be further verified by the tuple of overall cost and number of transmissions, which was observed to be $(7301.67, 499), (11683.64, 386), (16160.71, 374)$ and $(19002.01, 249)$ for increasing value of $\lambda$ in Fig. 4 over a horizon of 500 units. The corresponding signal trajectories for each case are provided in the Appendix.

### 6.1.4. Boeing flight control

Our final experimental setup is that of a 4-dimensional Boeing flight control system [28, 29]. We take 2-mode GMM with mean values $\begin{pmatrix} -5.0 & -4.0 & -3.0 & -2.0 \\ 4 & 5 & 3.0 & 2.0 \end{pmatrix}^{\top}$, covariance matrix $\begin{pmatrix} \text{diag}(0.1 & 0.2 & 0.3 & , 0.4) \\ \text{diag}(0.4 & 0.1 & 0.3 & 0.2) \end{pmatrix}$, and the weight vector as $(0.3, 0.7)$. The scatter plot demonstrating the scheduling landscape (for $\lambda = 45$) and the implicit scheduler-estimator agreement for the GMM modes, is presented in Fig. 5 along with its corresponding signal trajectories in Fig. 12 in the Appendix.
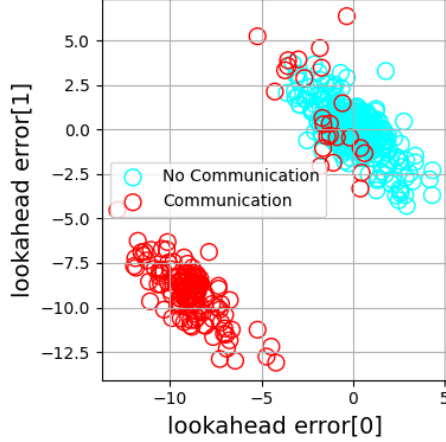
Figure 5: Scheduling landscape for the Boeing flight control system showing clear partitioning between communication and no-communication events.

## 6.2. Baseline Comparison

In this study, we fix the estimator to that trained using CALM, and compare the Pareto-fronts of the estimation-scheduling trade-offs produced using the following baselines (see Fig. 6).

- Periodic Scheduling: A fixed-period transmission scheme in which the plant communicates with the estimator at regular integer time intervals (in an open-loop manner). We evaluate period lengths of 1, 2 and 3, to observe trade-offs between communication frequency and estimation performance.

- Event-Triggered Policy: A threshold-based feedback scheduling policy of the form $\|e_t\|^2 \geq \tau$, where $\tau > 0$ is a design parameter. Communication occurs only when the estimation error exceeds the threshold.

From Fig. 6 (for $\lambda = 45$ on the tracking and pendulum problem), we observe that the trade-off performance of the proposed algorithm CALM (over two different sets of noise realizations) is better than the baselines on the periodic policy, and the event-triggered threshold policy.

## 7. Conclusion

We proposed a novel algorithmic framework for the joint learning of communication schedulers and state estimators in stochastic dynamical systems, using neural networks as function approximators. The approach, applicable to general nonlinear systems, alternates between training the estimator and the scheduler, the latter of which is optimized using proximal policy optimization. Our results highlight a key insight: communication silence—i.e., no-communication events—implicitly carries information that can be leveraged to enhance estimation accuracy, and close a long-standing gap in co-design methodologies. Extensive experiments across standard benchmarks, including inverted pendulum stabilization, Van
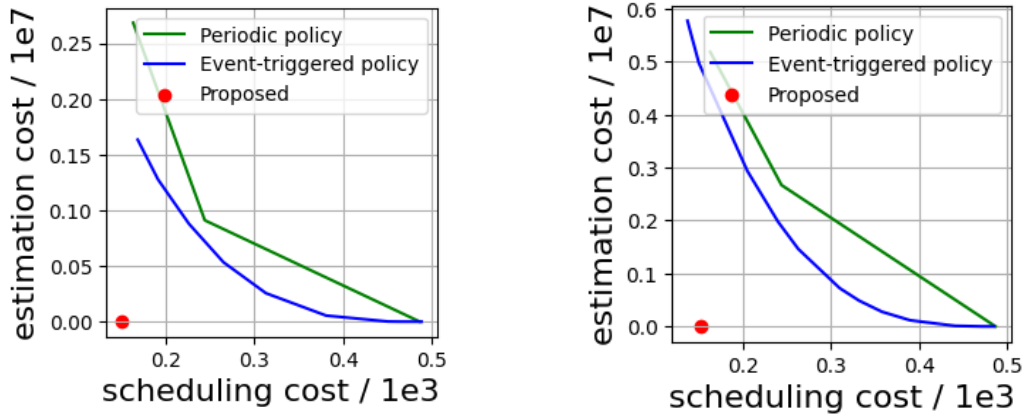
Figure 6: Estimation-scheduling trade-off for a fixed estimator for tracking (left) and inverted pendulum problem (right).

der Pol oscillator, and robot trajectory tracking systems, reveal that our method consistently outperforms classical estimators and widely used heuristic scheduling strategies such as periodic and event-triggered ones.

References

[1] Torsten Söderström. Discrete-Time Stochastic Systems: Estimation and Control. Springer Science & Business Media, 2012.

[2] Karaputugala Madushan Thilina, Kae Won Choi, Nazmus Saquib, and Ekram Hossain. Machine learning techniques for cooperative spectrum sensing in cognitive radio networks. IEEE Journal on Selected Areas in Communications, 31(11):2209–2221, 2013.

[3] Rudolph Emil Kalman. A new approach to linear filtering and prediction problems. 1960.

[4] Jakob Marschak. Elements for a theory of teams. Management Science, 1(2):127–137, 1955.

[5] Roy Radner. Team decision problems. The Annals of Mathematical Statistics, 33(3): 857–881, 1962.

[6] Feldbaum A.A. Dual Control Theory, pages 181–196. 2001. doi: 10.1109/9780470544334.ch10.

[7] Serdar Yüksel and Tamer Başar. Stochastic Networked Control Systems: Stabilization and Optimization under Information Constraints. Springer Science & Business Media, 2013.

[8] Orhan C Imer and Tamer Başar. Optimal estimation with limited measurements. International Journal of Systems, Control and Communications, 2(1-3):5–29, 2010.

[9] Gabriel M Lipsa and Nuno C Martins. Remote state estimation with communication costs for first-order LTI systems. IEEE Transactions on Automatic Control, 56(9): 2013–2025, 2011.

[10] Adam Molin. Optimal Event-Triggered Control with Communication Constraints. PhD thesis, Technische Universität München, 2014.

[11] Adam Molin and Sandra Hirche. Event-triggered state estimation: An iterative algorithm and optimality properties. IEEE Transactions on Automatic Control, 62(11): 5939–5946, 2017.

[12] Dipankar Maity and John S Baras. Minimal feedback optimal control of linear-quadratic-Gaussian systems: no communication is also a communication. IFAC-PapersOnLine, 53(2):2201–2207, 2020.

[13] Mark Eisen, Santosh Shukla, Dave Cavalcanti, and Amit S Baxi. Communication-control co-design in wireless edge industrial systems. In 2022 IEEE 18th International Conference on Factory Communication Systems (WFCS), pages 1–8. IEEE, 2022.

[14] Siyi Wang and Sandra Hirche. Infinite-horizon optimal scheduling for feedback control. arXiv preprint arXiv:2402.08819, 2024.

[15] Chithrupa Ramesh, Henrik Sandberg, and Karl H Johansson. Design of state-based schedulers for a network of control loops. IEEE Transactions on Automatic Control, 58(8):1962–1975, 2013.

[16] Michael Hertneck, David Meister, and Frank Allgöwer. Current trends and future directions in event-based control. arXiv preprint arXiv:2505.22378, 2025.

[17] Touraj Soleymani, John S Baras, Sandra Hirche, and Karl H Johansson. Value of information in feedback control: Global optimality. IEEE Transactions on Automatic Control, 68(6):3641–3647, 2022.

[18] Shubham Aggarwal, Dipankar Maity, and Tamer Başar. InterQ: A DQN framework for optimal intermittent control. IEEE Control Systems Letters, 2025.

[19] Ya-Chien Chang, Nima Roohi, and Sicun Gao. Neural Lyapunov control. Advances in Neural Information Processing Systems, 32, 2019.

[20] Ruikun Zhou, Thanin Quartz, Hans De Sterck, and Jun Liu. Neural Lyapunov control of unknown nonlinear systems with stability guarantees. Advances in Neural Information Processing Systems, 35:29113–29125, 2022.

[21] Junlin Wu, Andrew Clark, Yiannis Kantaros, and Yevgeniy Vorobeychik. Neural Lyapunov control for discrete-time systems. Advances in Neural Information Processing Systems, 36:2939–2955, 2023.

[22] Touraj Soleymani, John S Baras, and Sandra Hirche. Value of information in feedback control: Quantification. IEEE Transactions on Automatic Control, 67(7):3730–3737, 2021.

[23] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347, 2017.

[24] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In International Conference on Machine Learning, pages 1889–1897. PMLR, 2015.

[25] Xiaodong Zhang and Nikolay Dimitrov. Extreme wind turbine response extrapolation with the gaussian mixture model. Wind Energy Science, 8(10):1613–1623, 2023.

[26] Vanessa Hernandez-Cruz, Xiaotong Zhang, and Kamal Youcef-Toumi. Bayesian intention for enhanced human robot collaboration. arXiv preprint arXiv:2410.00302, 2024.

[27] John L Iannamorelli and Keith A LeGrand. Adaptive gaussian mixture filtering for multi-sensor maneuvering cislunar space object tracking. The Journal of the Astronautical Sciences, 72(1):2, 2025.

[28] Stephen Boyd and Lieven Vandenberghe. Introduction to Applied Linear Algebra: Vectors, Matrices, and Least Squares. Cambridge University Press, 2018.

[29] Shahab Ataei, Dipankar Maity, and Debdipta Goswami. Qsid-mpc: Model predictive control with system identification from quantized data. arXiv preprint arXiv:2503.19102, 2025.

[30] Hassan K Khalil. Nonlinear Systems, volume 3. Prentice Hall, Upper Saddle River, NJ, 2002.

## Appendix: Supplementary Material

## 8. Further Details about Experiments

### 8.1. System Dynamics

We describe all the system dynamics in continuous-time. Subsequently, we discretize them by setting the discretization interval to 0.05 seconds in all experiments.

#### 8.1.1. Inverted Pendulum

We use the following stochastic discrete-time dynamics model for the inverted pendulum system [19, 20, 21]:

$$x_{t+1} = \begin{pmatrix} 1 & \varepsilon \\ \frac{g}{\ell}\varepsilon & 1 - \frac{b}{m\ell^2}\varepsilon \end{pmatrix} x_t + \begin{pmatrix} 0 \\ \frac{1}{m\ell^2}\varepsilon \end{pmatrix} u_t + w_t,$$

where the state $x := [\theta, \dot{\theta}]^\top$ captures the angle $\theta$ and the angular velocity $\dot{\theta}$. The parameter $g$ denotes gravity, $m$ denotes mass of the ball, $b$ denotes friction coefficient, and $\ell$ denotes pendulum length. For all experiments, we take $g = 9.81$, $m = 0.15$, $b = 0.1$, $\ell = 0.5$, and $\varepsilon = 0.05$.

### 8.1.2. Van der Pol (VdP) oscillator

We use the following nonlinear dynamics for the VdP oscillator [30]:

$$(x_1)_{t+1} = (x_1)_t + \varepsilon[(x_2)_t + (w_1)_t]$$
$$(x_2)_{t+1} = (x_2)_t + \varepsilon[\mu(1 - (x_1)_t^2)(x_2)_t - (x_2)_t + (w_2)_t]$$

with $x_1$ and $x_2$ denoting the position and velocity of the second-order system with the state $x = [x_1,\ x_2]^\top$. Further, $\mu > 0$ denotes the damping strength. For all experiments, we take the parameters as $\mu = 0.025$ and $\varepsilon = 0.05$.

### 8.1.3. Robot trajectory tracking system

We use the following discretized model for trajectory-tracking system [21]:

$$x_{t+1} = \begin{pmatrix} 1 & 2 \\ 0 & 1 - 0.04\varepsilon \end{pmatrix} x_t + \begin{pmatrix} 0 \\ \varepsilon \end{pmatrix} u_t + w_t,$$

where the state $x$ constitutes the distance and the angle errors from the corresponding reference values, the control input constitutes the steering angle $\delta$. For all experiments, we take $\varepsilon = 0.05$.

### 8.1.4. Boeing 747 flight control

We use the following 4-dimensional longitudinal flight control system of the Boeing 747 aircraft in steady level flight at an altitude of 40,000 ft, a speed 774 ft/s, and a time unit of one second [28, 29]:

$$x_{t+1} = \begin{pmatrix} 0.99 & 0.03 & -0.02 & -0.32 \\ 0.01 & 0.47 & 4.7 & 0.0 \\ 0.02 & -0.06 & 0.40 & 0.0 \\ 0.01 & -0.04 & 0.72 & 0.99 \end{pmatrix} x_t + \begin{pmatrix} 0.01 & 0.99 \\ -3.44 & 1.66 \\ -0.83 & 0.44 \\ -0.47 & 0.25 \end{pmatrix} u_t + w_t$$

The state $x_t$ constitutes deviations (from the nominal values) of the velocity along the aircraft body axis, velocity perpendicular to the body axis, angle of the body above horizontal, and derivative of the angle of the body of the body axis. The control input constitutes the deviations of the elevator angle and the engine thrust, from the nominal values.

### 8.2. Controlled systems

For systems with explicit control inputs such as the inverted pendulum, trajectory tracking, and the flight control ones, we apply the infinite-horizon linear-quadratic state feedback control solution. Precisely, for given matrices $Q \succeq 0$ and $R \succ 0$, we apply the control input:

$$u_t = Kx_t$$

at the dynamical system while at the estimator, which has access to only the best estimate of the state $\hat{x}_t$ at time $t$, we apply the input:

$$u_t = K\hat{x}_t.$$

Here, we have

$$K = -\gamma(P + \gamma B^\top RB)^{-1}B^\top PA$$

and $P \succeq 0$ is the unique positive semi-definite solution to the algebraic Riccati equation:

$$P = \gamma A^\top PA + Q - \gamma^2 A^\top PB(B^\top RB + P)^{-1}B^\top PA.$$

For all experiments, we used both Q and R to be identity matrices of suitable dimensions and the discount factor to be 0.9999.

### 8.3. Training Hyperparameters

The training parameters for CALM were set to the following:
   Estimator NN activation function: ReLU
   PPO network activation function: ReLU
   Optimizer for Estimator NN: Adam with weight decay
   Optimizer for PPO: Adam
   learning rate for both networks: 1e-3
   GAE lambda = 0.9
   PPO clipping parameter = 0.2
   Number of outer epochs: 10
   Number of inner epochs for PPO: 80
   Number of inner epochs for estimator: 150
   Number of epochs for scheduler policy training (PPO) with fixed linear estimator: 1000
   Trajectory horizon length: 80

## 9. Extended Experiments

### 9.1. Inverted Pendulum

In Fig. 7, we plot the state, state estimate and the estimation error trajectories (overlayed with the communication instances) corresponding to the (left) scatter plot in Fig. 2 of a 2-mode GMM noise distribution. Additionally, in Fig. 8, we also plot a comparative scatter plot obtained by training and evaluating a baseline linear estimator on a 2- and a 3-mode GMM noise distribution similar to the CALM algorithm results of Fig. 2. The corresponding signal trajectories for the 2-mode GMM case are also plotted in FIg. 9.

### 9.2. Van der Pol

In Fig. 10, we plot the state, state estimate and the estimation error trajectories, overlapped with the communication instances (red markers) corresponding to the scatter plot of the left subfigure in Fig. 3.

### 9.3. Robot trajectory tracking

For the trajectory tracking system, we plot the signal trajectories for each value of $\lambda$ corresponding to Fig. 11. As $\lambda$ increases from 15 to 30 to 40 to 70, the number of communication instances decrease from 499 to 386 to 374 to 249, as aligned with intuition.
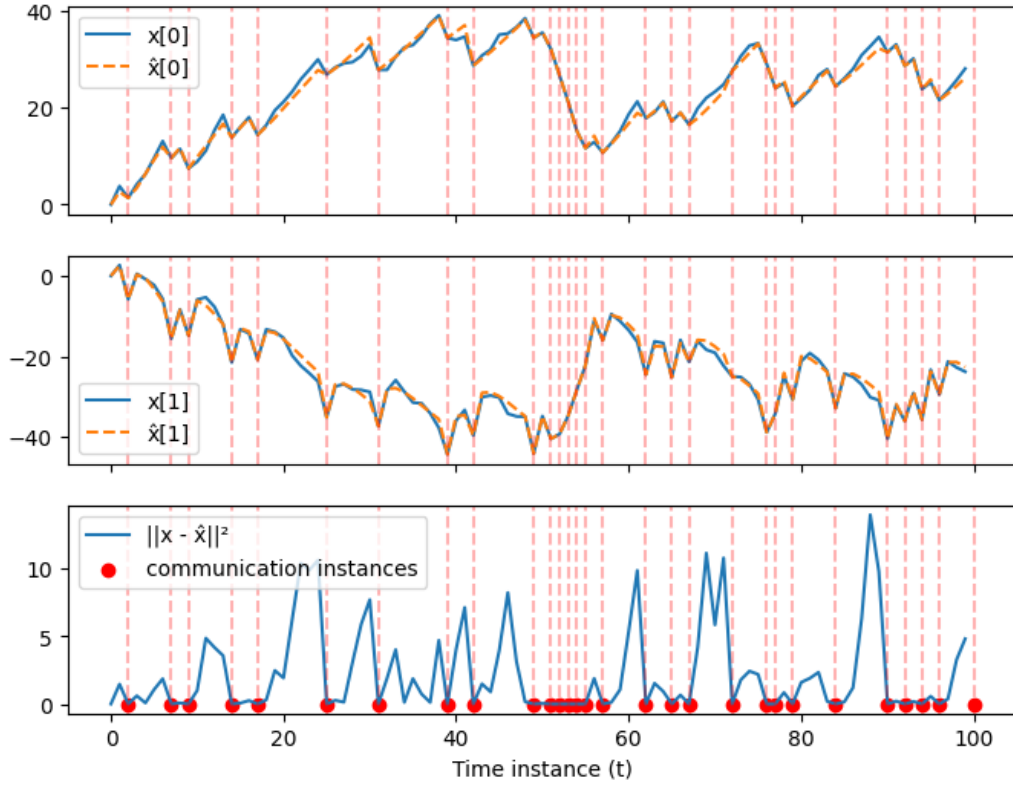
Figure 7: Signal trajectories with communication instances (red markers) corresponding to the left subfigure in Fig. 2 for the inverted pendulum.
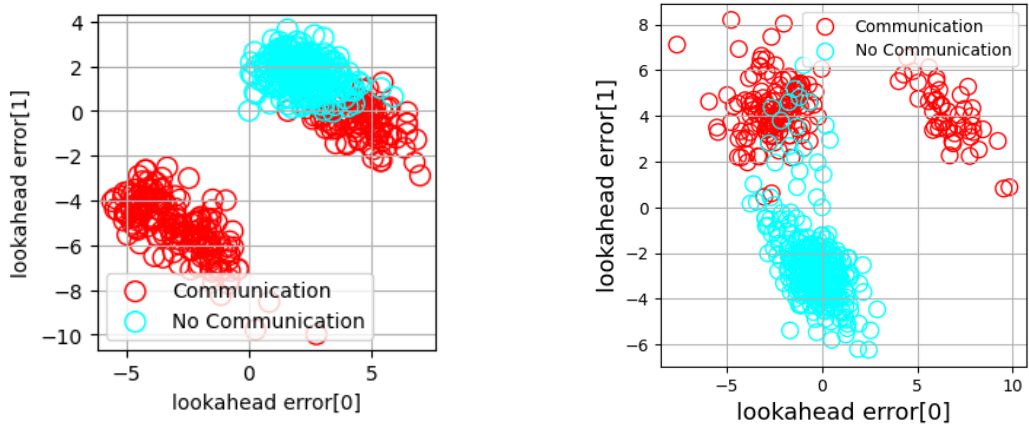


Figure 8: Scheduling landscapes for a 2-mode GMM (left) and 3-mode GMM (right) for the inverted pendulum system (red denotes transmissions while cyan denotes silence).
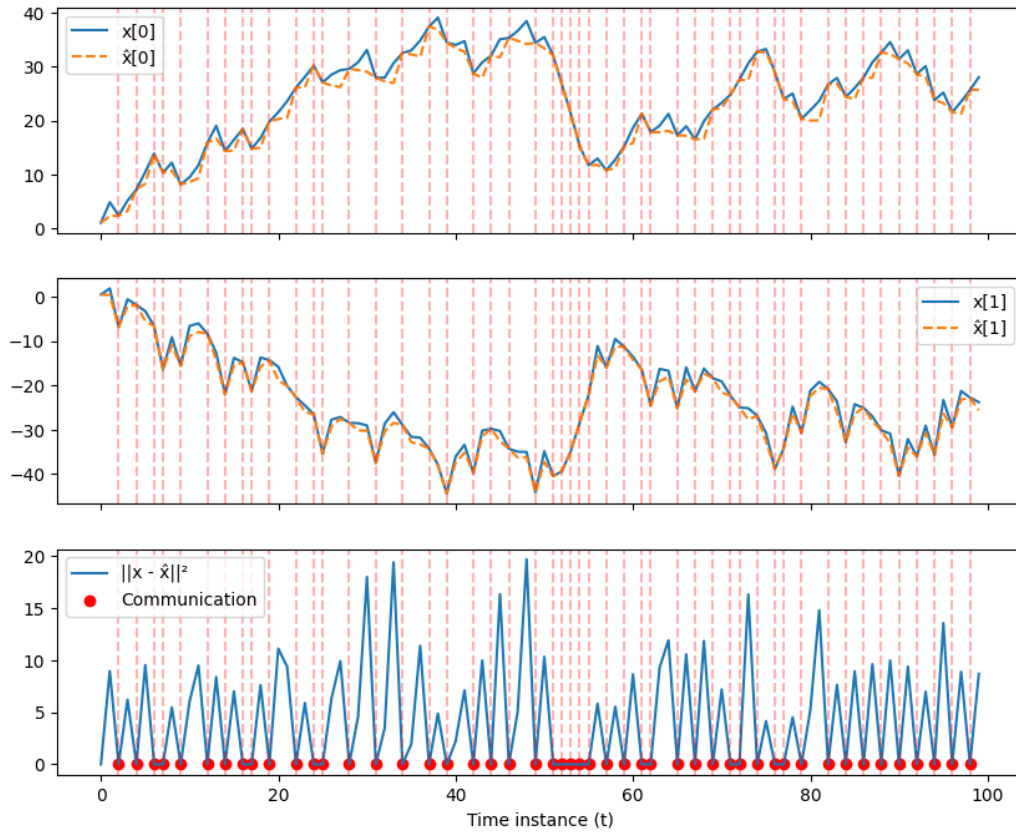
Figure 9: Signal trajectories with communication instances (red markers) corresponding to (the left subfigure in) Fig. 8 for the inverted pendulum system.
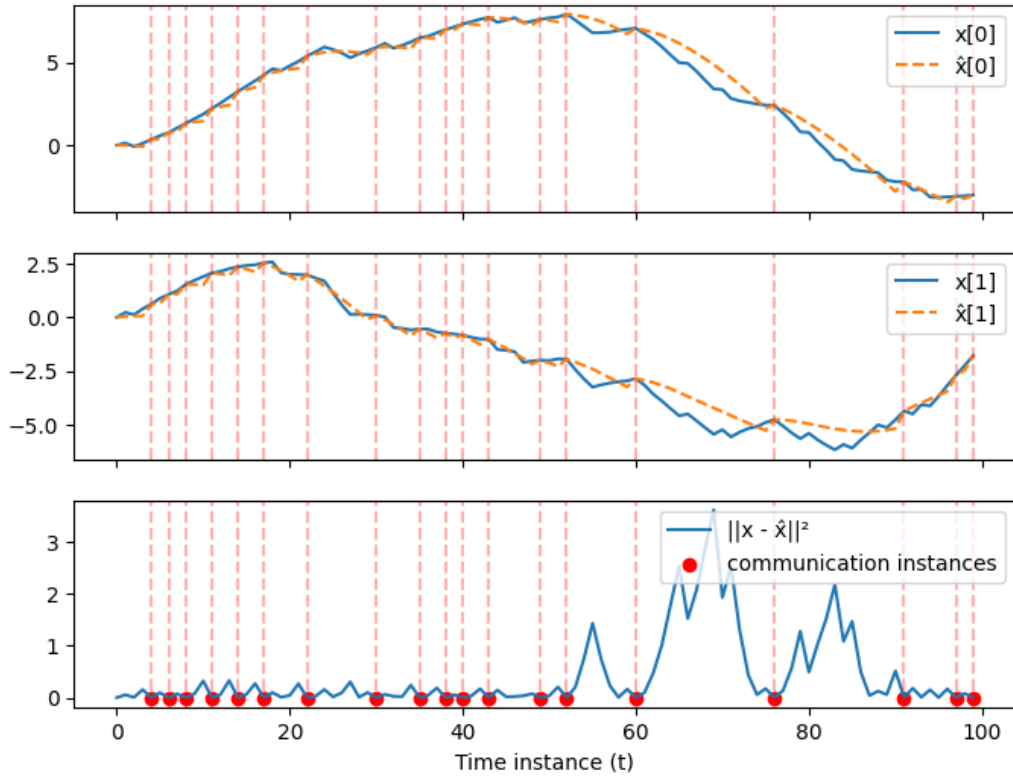
Figure 10: Signal trajectories with communication instances (red markers) corresponding to the left subfigure in Fig. 3 ($\lambda = 0.7$) for the VdP oscillator.
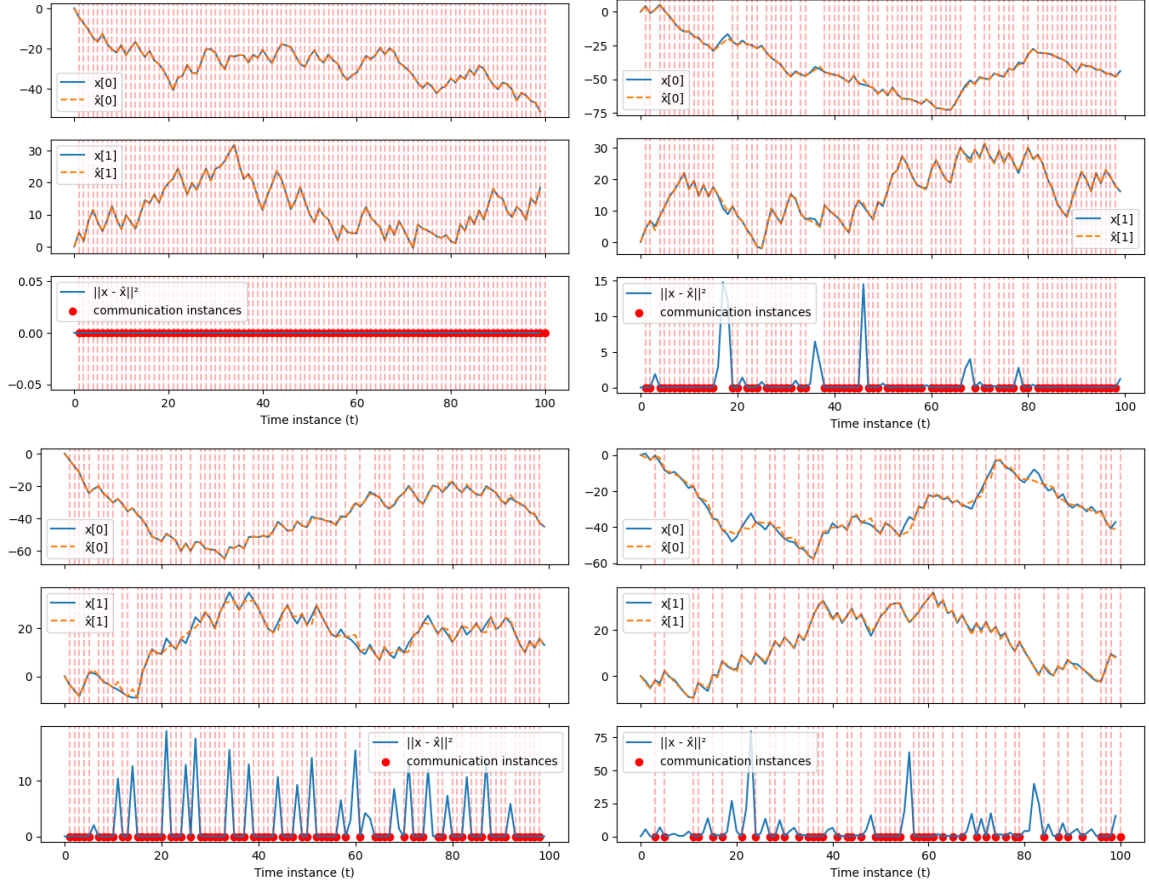
Figure 11: Signal trajectories over time for the robot trajectory tracking system: $\lambda = 15$ (top left), $\lambda = 30$ (top right), $\lambda = 40$ (bottom-left) and $\lambda = 70$ (bottom right).
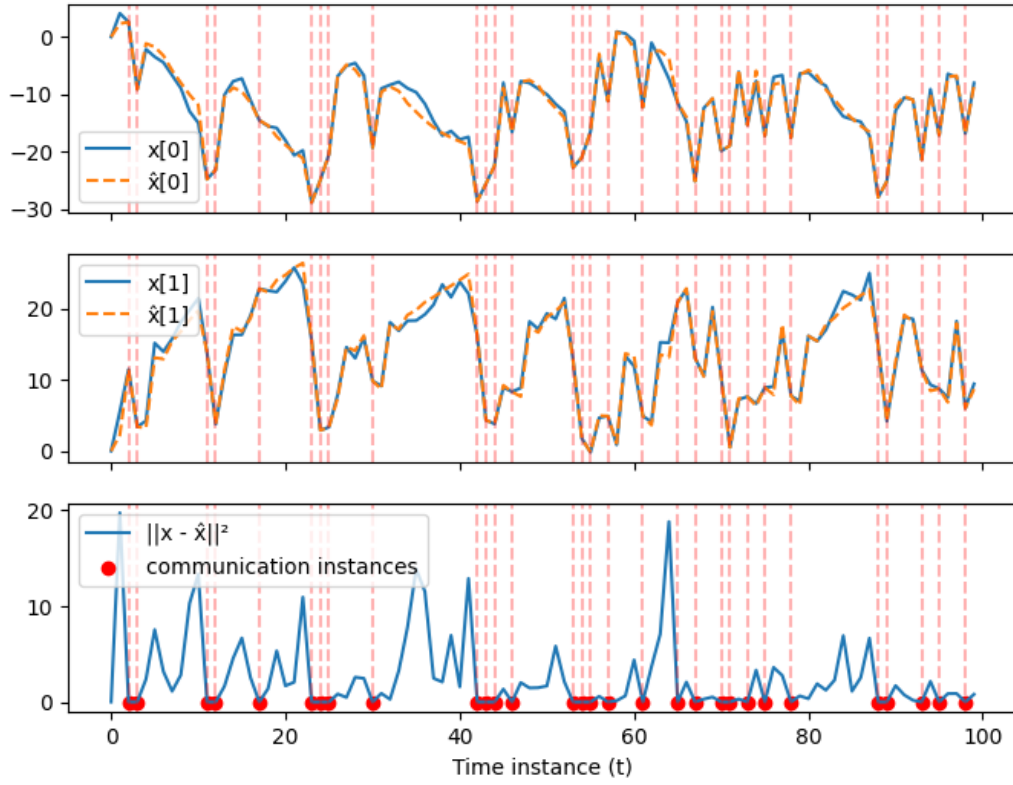
Figure 12: Signal trajectories with communication instances (red markers) corresponding to Fig. 5 for the Boeing flight control system.