

Entropy-Regularized Partially Observed Markov Decision Processes

Timothy L. Molloy, *Member, IEEE*, and Girish N. Nair, *Fellow, IEEE*

Abstract

We investigate partially observed Markov decision processes (POMDPs) with cost functions regularized by entropy terms describing state, observation, and control uncertainty. Standard POMDP techniques are shown to offer bounded-error solutions to these entropy-regularized POMDPs, with exact solutions possible when the regularization involves the joint entropy of the state, observation, and control trajectories. Our joint-entropy result is particularly surprising since it constitutes a novel, tractable formulation of active state estimation.

I. INTRODUCTION

Partially observed Markov decision processes (POMDPs) and Markov decision processes (MDPs) with information-theoretic costs have attracted widespread attention across systems and control [2]–[5], computer science [6]–[8], signal processing [9]–[12], and robotics [13]–[15]. Interest in such POMDPs has been driven, in large part, by *active state estimation* problems in which information-theoretic costs describing the uncertainty about latent states are minimized in order to aid or enhance the performance of state estimation algorithms [5], [6], [9], [10]. Interest in such MDPs has, in contrast, been driven by a desire within applications such as networked control and economics to develop control policies (or decision-makers) that are *rationally inattentive* or “data-frugal” in that they trade-off control performance to reduce the (data) rate at which state information is used to make control decisions [2], [3], [8]. Despite interest in rate-cost trade-offs in MDPs, limited attention has been paid to similar problems in POMDPs.

The first author was with the Dept. of Electrical and Electronic Engineering, University of Melbourne, VIC 3010, Australia. He is now with the CIICADA Lab, School of Engineering, The Australian National University (ANU), Canberra, ACT 0200, Australia (e-mail: timothy.molloy@anu.edu.au) The second author is with the Dept. of Electrical and Electronic Engineering, University of Melbourne, VIC 3010, Australia (e-mail: gnair@unimelb.edu.au)

This work received funding from the Australian Government, via grant AUSMURIB000001 associated with ONR MURI grant N00014-19-1-2571.

Preliminary versions of some results in this paper were presented at the 2022 American Control Conference [1].

Motivated by data-frugal POMDPs with potential applications to active state estimation, we investigate POMDPs with information-theoretic entropy costs that penalize observation incompressibility and/or state uncertainty.

POMDPs with information-theoretic costs have been extensively investigated for active state estimation, with popular costs including the (negative) mutual information between states and observations [4], the (Shannon or Rényi) entropy of Bayesian filter estimates [6], [9], [13], and the entropy of Bayesian smoother or Viterbi algorithm estimates [5], [14], [15] (see [10, Chapter 8] and references therein for more). These POMDPs have been shown to be amenable to bounded-error (approximate) solution using standard POMDP solvers when their cost and cost-to-go (or value functions) are concave and admit piecewise-linear concave (PWLC) approximations (cf. [10, Section 8.4.4], [6]). For example, we recently showed that the *smoother entropy* (i.e., the conditional entropy of the state trajectory given observations and controls) can be (approximately) optimized in this manner [5], [16].

Outside of applications involving active state estimation, POMDPs with information-theoretic costs have received only limited specialized attention. Most notably, in Bayesian experimental design (involving degenerate POMDPs with a time-invariant or constant state), the entropy of the observations has been explored as a cost to encourage the selection of controls (i.e. experiments) with predictable outcomes (see [17], [18]). In the context of linear-quadratic regulators and linear-quadratic-Gaussian control (i.e. POMDPs with specialized linear dynamics but continuous state, control, and observation spaces), the *directed information* from the observations to controls has been used as a cost to study the trade-off between feedback (data) rates and control costs (see [19]–[21]). Similarly, in the context of MDPs (i.e. degenerate POMDPs with fully observed states), various information-theoretic quantities such as the directed information, mutual information, and *transfer entropy*, have been used as costs to penalize feedback from the states to the controls so as to study rate-cost trade-offs [2], [3], [8]. Despite growing interest, solving POMDPs (and MDPs) with rate-cost trade-offs has proved difficult due to complications including randomized policies [2], [17]–[19], the design of observation processes [2], [3], [8], [19], and the need to solve nonconvex optimization problems [2], [7].

The main contribution of this paper is the proposal of POMDPs with costs regularized by combinations of the *input-output entropy* (i.e. the entropy of the observations and controls) and the smoother entropy. Such *entropy-regularized* POMDPs (ERPOMDPs) are novel in that they both introduce rate-cost trade-offs into standard POMDPs (due to the relevance of the input-output entropy to rate-cost trade-offs via Shannon’s source coding theorem), and generalize recent work on active state estimation involving the smoother entropy to include compressible (or predictable) observations (cf. [5], [16]). Importantly, we show that ERPOMDPs admit bounded-error PWLC solutions via standard POMDP techniques, in general,

and *exact* solutions *without* additional PWLC approximations in the special case where the smoother and input-output entropies are equally weighted and become the joint entropy of the states, observations, and controls. The solution of ERPOMDPs involving the joint entropy without PWLC approximations is surprising since the vast majority of other POMDPs with information-theoretic costs are entirely intractable without them. Compared to our preliminary work in [1], significant extensions in this paper include: 1) consideration of a generalized problem (7) with arbitrary combinations of the smoother and input-output entropies; 2) development of entirely new results in Lemmas 3.1 and 3.2, and Theorem 3.1 concerning the input-output entropy; and, 3) new operational interpretations of ERPOMDPs in Section V.

This paper is structured as follows. In Section II, we pose ERPOMDPs and examine their solution via standard POMDP techniques in Section III. In Section IV, we examine *exact* ERPOMDP solutions in the case of joint-entropy regularization. Finally, we provide interpretations of ERPOMDPs in Section V, simulations in Section VI, and conclusions in Section VII.

Notation: Random variables will be denoted by capital letters (e.g., X), their realizations by lower case letters (e.g., x), and associated sequences by letters with superscripts denoting their final times (e.g., $X^T \triangleq \{X_0, X_1, \dots, X_T\}$ and $x^T \triangleq \{x_0, x_1, \dots, x_T\}$). The probability mass function (pmf) of X will be written $p(x)$, the joint pmf of X and Y written $p(x, y)$, and the conditional pmf of X given $Y = y$ written $p(x|y)$ or $p(x|Y = y)$. For a function f of X , the expectation of f is $E_X[f(X)]$ and the conditional expectation of f under $p(x|y)$ is $E[f(X)|y]$. The *pointwise* conditional entropy of X given y is $H(X|y) \triangleq -E[\log p(X|Y = y)|y]$, and the conditional entropy of X given Y is $H(X|Y) \triangleq E_Y[H(X|y)]$, with the base of the logarithm being 2.

II. PROBLEM FORMULATION

Let X_k for $k \geq 0$ be a discrete-time first-order Markov chain with the finite state space $\mathcal{X} \triangleq \{1, 2, \dots, N_x\}$. Let the initial state X_0 be distributed according to the pmf $\rho \in \Delta$ with components $\rho(x_0) \triangleq P(X_0 = x_0)$ where $\Delta \triangleq \{\rho \in [0, 1]^{N_x} : \sum_{x \in \mathcal{X}} \rho(x) = 1\}$ is the $(N_x - 1)$ -dimensional probability simplex. Let the (controlled) transition dynamics of the state X_k be described by the transition kernel:

$$A^{x, \bar{x}}(u) \triangleq p(X_{k+1} = x | X_k = \bar{x}, U_k = u) \quad (1)$$

for $k \geq 0$ with the controls $U_k = u$ belonging to the finite set $\mathcal{U} \triangleq \{1, 2, \dots, N_u\}$. The state process X_k is (partially) observed through a stochastic observation process Y_k for $k \geq 0$ taking values in the finite

set $\mathcal{Y} \triangleq \{1, 2, \dots, N_y\}$. The observations Y_k are distributed according to the kernel:

$$B^{x,y}(u) \triangleq p(Y_k = y | X_k = x, U_{k-1} = u) \quad (2)$$

for $k > 0$ with $B^{x_0,y_0} \triangleq p(Y_0 = y_0 | X_0 = x_0)$. The controls U_k for $k \geq 0$ arise from a potentially stochastic output-feedback policy $\mu \triangleq \{\mu_k^{i_k} : k \geq 0\}$ with (conditional) pmfs

$$\mu_k^{i_k}(u_k) \triangleq p(U_k = u_k | Y^k = y^k, U^{k-1} = u^{k-1})$$

where $i_k \triangleq (y^k, u^{k-1})$ is a realization of the *information state* $I_k \triangleq (Y^k, U^{k-1})$. The joint pmf of (X^T, I_T) under μ is

$$\begin{aligned} p_\mu(x^T, y^T, u^{T-1}) &= \rho(x_0) B^{x_0, y_0} \\ &\times \prod_{k=0}^{T-1} A^{x_{k+1}, x_k}(u_k) \mu_k^{i_k}(u_k) B^{x_{k+1}, y_{k+1}}(u_k), \end{aligned} \quad (3)$$

for $T > 0$ where $\prod_{k=0}^{-1}$ is taken as the identity matrix. We denote expectation under p_μ as $E_\mu[\cdot]$. A policy $\mu = \{\mu_k^{i_k} : k \geq 0\}$ is *deterministic* if, at all times $k \geq 0$, the support of $\mu_k^{i_k}$ is concentrated at a single control u_k ; otherwise μ is *stochastic*. Let the set of all policies (stochastic or deterministic) be \mathcal{P} .

To introduce our ERPOMDP problem, let us define the smoother entropy for $T \geq 0$ under a policy $\mu \in \mathcal{P}$ as

$$H_\mu(X^T | Y^T, U^{T-1}) \triangleq -E_\mu[\log p_\mu(X^T | Y^T, U^{T-1})], \quad (4)$$

and let us define the input-output entropy under μ as

$$H_\mu(Y^T, U^{T-1}) \triangleq -E_\mu[\log p_\mu(Y^T, U^{T-1})]. \quad (5)$$

Let us also define the additive cost functional

$$J_\mu^T \triangleq E_\mu \left[c_T(X_T) + \sum_{k=0}^{T-1} c(X_k, U_k) \right] \quad (6)$$

where $c : \mathcal{X} \times \mathcal{U} \mapsto \mathbb{R}$ and $c_T : \mathcal{X} \mapsto \mathbb{R}$ are arbitrary cost functions dependent on the state and control values.

Our ERPOMDP problem is to find a policy that solves

$$\begin{aligned} \inf_{\mu \in \mathcal{P}} \quad & E_T[J_\mu^T + \beta H_\mu(X^T | Y^T, U^{T-1}) + \lambda H_\mu(Y^T, U^{T-1})] \\ \text{s.t.} \quad & X_{k+1} | X_k, U_k \sim A^{x_{k+1}, x_k}(u_k), \quad X_0 \sim \rho \\ & Y_{k+1} | X_{k+1}, U_k \sim B^{x_{k+1}, y_{k+1}}(u_k), \quad Y_0 | X_0 \sim B^{x_0, y_0} \\ & U_k | I_k \sim \mu_k^{i_k}(u_k) \end{aligned} \quad (7)$$

for given nonnegative constants $\beta, \lambda \geq 0$ where the horizon $T \geq 0$ is a random variable with a geometric distribution with (probability of nonoccurrence) parameter $0 < \gamma < 1$ such that $\zeta_t \triangleq P(T = t) = \gamma^t(1-\gamma)$ for $t \geq 0$. Despite the entropies above not being in additive forms, we shall show later that the total cost over a geometrically distributed finite horizon is equivalent to a discounted additive cost over an infinite horizon, with γ being the discount factor (cf. [22]).

The motivation behind our ERPOMDP problem (7) is twofold. Firstly, the input-output entropy has an interpretation as the minimum expected number of bits required to transmit or store the observations and controls (Y^T, U^{T-1}) (via Shannon's source coding theorem [23, Section 5.5]). Solving (7) with $\lambda > 0$ (and $\beta \geq 0$) thus leads to policies that reduce the number of bits used for feedback control (similar to the MDPs in [2]). Secondly, the smoother entropy intuitively describes the uncertainty associated with estimates of the states X^T given the observations and controls. Solving (7) with $\beta > 0$ and $\lambda > 0$ thus leads to policies that reduce the number of bits used to store the observations and controls, whilst ensuring that the states can still be estimated from them. Operational interpretations of (7) are discussed further in Section V.

Solving (7) is greatly simplified if we are able to use standard POMDP solution techniques since they are increasingly able to handle large-scale problems (cf. [6], [7], [24], [25]). As discussed in [6] and [10, Chapter 8], the use of standard POMDP techniques to find bounded-error solutions to (7) requires that: 1) its cost function can be written as an additive function of a sufficient statistic of the information state known as the *belief state*; 2) it can be reformulated as a (fully observed) MDP in terms of the belief state with cost functions that can be arbitrarily well-approximated by PWLC functions; and 3) it can be solved by deterministic policies. In [5], [16], we showed that this solution approach is possible without the input-output entropy (i.e. when $\beta > 0 = \lambda$) by establishing a belief-state expression of the smoother entropy. The input-output entropy appears more challenging to optimize since its naive factorization as $H_\mu(Y^T, U^{T-1}) = H_\mu(Y^T | U^{T-1}) + H_\mu(U^{T-1})$ shows immediately that it involves the (unconditional) entropy of the policies $H_\mu(U^{T-1})$, which means that we must consider the possibility of optimal policies solving (7) being stochastic. We shall therefore focus on: 1) establishing a belief-state expression of the input-output entropy; 2) showing that it suffices to consider deterministic policies in solving (7); and, 3) developing belief MDP reformulations of (7).

III. BELIEF-STATE FORMS AND MDP REFORMULATION

In this section, we revisit the concept of the belief state and a belief-state form of the smoother entropy. We then establish a novel belief-state form of the input-output entropy that enables (7) to be reformulated as a belief MDP amenable to bounded-error solution using standard POMDP techniques.

A. Belief State and Smoother Entropy

Let $\pi_k \in \Delta$ with $\pi_k(x) \triangleq p(X_k = x|y^k, u^{k-1})$ for $x \in \mathcal{X}$ be the belief state, which evolves via the Bayesian filter:

$$\pi_{k+1}(x) = \frac{B^{x, y_{k+1}}(u_k) \sum_{\bar{x} \in \mathcal{X}} A^{x, \bar{x}}(u_k) \pi_k(\bar{x})}{\sum_{\underline{x} \in \mathcal{X}} \sum_{\bar{x} \in \mathcal{X}} B^{\underline{x}, y_{k+1}}(u_k) A^{\underline{x}, \bar{x}}(u_k) \pi_k(\bar{x})}$$

for $k \geq 0$ with $\pi_0(x) = B^{x, y_0} \rho(x) / (\sum_{\bar{x} \in \mathcal{X}} B^{\bar{x}, y_0} \rho(\bar{x}))$ for $x \in \mathcal{X}$. We write the filter as $\pi_{k+1} = \Pi(\pi_k, u_k, y_{k+1})$.

In [5], [16], we showed that the smoother entropy satisfies

$$H_\mu(X^T|Y^T, U^{T-1}) = E_\mu \left[\tilde{G}_1(\pi_T) + \sum_{k=0}^{T-1} \tilde{G}_2(\pi_k, U_k) \right] \quad (8)$$

where $\tilde{G}_1(\pi_k) \triangleq -\sum_{x \in \mathcal{X}} \pi_k(x) \log \pi_k(x)$ is the *belief-state entropy*, i.e. $H_\mu(X_k|y^k, u^{k-1})$, and

$$\begin{aligned} \tilde{G}_2(\pi_k, u_k) &\triangleq \sum_{x, \bar{x} \in \mathcal{X}} A^{x, \bar{x}}(u_k) \pi_k(\bar{x}) \log \sum_{\underline{x} \in \mathcal{X}} \frac{A^{x, \underline{x}}(u_k) \pi_k(\underline{x})}{A^{x, \bar{x}}(u_k) \pi_k(\bar{x})} \\ &= H_\mu(X_k, X_{k+1}|y^k, u^k) - H_\mu(X_{k+1}|y^k, u^k) \end{aligned}$$

is the difference between the entropy of $p_\mu(x_k, x_{k+1}|y^k, u^k) = A^{x_{k+1}, x_k}(u_k) \pi_k(x_k)$, i.e. $H_\mu(X_k, X_{k+1}|y^k, u^k)$, and the entropy of $p_\mu(x_{k+1}|y^k, u^k) = \sum_{x_k \in \mathcal{X}} A^{x_{k+1}, x_k}(u_k) \pi_k(x_k)$, i.e. $H_\mu(X_{k+1}|y^k, u^k)$, with these pmfs computed in the prediction step of the Bayesian filter. Specifically, the belief-state expression (8) arises because the pmf $p_\mu(x^T|y^T, u^{T-1})$ in (4) factorizes as

$$p_\mu(x^T|y^T, u^{T-1}) = \prod_{k=0}^{T-1} p_\mu(x_k|x_{k+1}^T, y^T, u^{T-1}) \quad (9)$$

via the chain rule with $x_{k+1}^T \triangleq \{x_{k+1}, \dots, x_T\}$, $x_{T+1}^T \triangleq \emptyset$, and since $p_\mu(x_k|x_{k+1}^T, y^T, u^{T-1}) = p_\mu(x_k|x_{k+1}, y^k, u^k) = p_\mu(x_k, x_{k+1}|y^k, u^k) / p_\mu(x_{k+1}|y^k, u^k)$ via the Markov property of the state and the structure of the measurement kernel and control policy. To reformulate (7) as a belief MDP, we need a similar expression for the input-output entropy.

B. Belief-State Form of Input-Output Entropy

To establish a novel belief-state form of the input-output entropy (5), we employ *causally conditioned entropies* as introduced by Kramer [26]. Let the causally conditioned entropy of Y^T given U^{T-1} under any policy $\mu \in \mathcal{P}$ be

$$H_\mu(Y^T||U^{T-1}) \triangleq \sum_{k=0}^{T-1} H_\mu(Y_k|Y^{k-1}, U^{k-1}) \quad (10)$$

where $H_\mu(Y_0|Y^{-1}, U^{-1}) \triangleq H(Y_0)$ is independent of μ . Similarly, let the causally conditioned entropy of U^{T-1} given Y^{T-1} under any policy $\mu \in \mathcal{P}$ be

$$H_\mu(U^{T-1}||Y^{T-1}) \triangleq \sum_{k=0}^{T-1} H_\mu(U_k|U^{k-1}, Y^k) \quad (11)$$

with $H_\mu(U_0|U^{-1}, Y^0) \triangleq H_\mu(U_0|Y_0)$, $H_\mu(U^{-1}||Y^{-1}) \triangleq 0$. Intuitively, $H_\mu(Y^T||U^{T-1})$ describes the uncertainty associated with the observations given the information causally gained from the controls, whilst $H_\mu(U^{T-1}||Y^{T-1})$ describes the uncertainty associated with the controls given the information causally gained from the (past) observations. The following lemma shows that the input-output entropy is the sum of these two causally conditioned entropies.

Lemma 3.1: For any $\mu \in \mathcal{P}$ and $T \geq 0$, we have that:

$$H_\mu(Y^T, U^{T-1}) = H_\mu(Y^T||U^{T-1}) + H_\mu(U^{T-1}||Y^{T-1}).$$

Proof: The proof is via induction. Note first that $H_\mu(Y^0||U^{-1}) + H_\mu(U^{-1}||Y^{-1}) = H(Y_0)$ proving the lemma assertion for $T = 0$. Assuming that the lemma assertion holds for trajectories shorter than some length $T > 0$, we now consider it for T . From (10) and (11),

$$\begin{aligned} & H_\mu(Y^T||U^{T-1}) + H_\mu(U^{T-1}||Y^{T-1}) \\ &= H_\mu(Y^{T-1}||U^{T-2}) + H_\mu(U^{T-2}||Y^{T-2}) \\ &\quad + H_\mu(Y_T|Y^{T-1}, U^{T-1}) + H_\mu(U_{T-1}|U^{T-2}, Y^{T-1}) \\ &= H_\mu(Y^{T-1}, U^{T-2}) + H_\mu(Y_T|Y^{T-1}, U^{T-1}) \\ &\quad + H_\mu(U_{T-1}|U^{T-2}, Y^{T-1}) = H_\mu(Y^T, U^{T-1}) \end{aligned}$$

where the second equality holds via the induction hypothesis, and the last equality holds due to the chain rule for conditional entropy. The proof of the lemma via induction is complete. \blacksquare

Lemma 3.1 differs from trivial expressions of the input-output entropy such as $H_\mu(Y^T, U^{T-1}) = H_\mu(Y^T) + H_\mu(U^{T-1}|Y^T) = H_\mu(U^{T-1}) + H_\mu(Y^T|U^{T-1})$ since these involve conditional entropies conditioned on the entire trajectories Y^T and U^{T-1} , whilst Lemma 3.1 establishes a form involving sums of conditional entropies only conditioned on the information state I_k at each time k . Lemma 3.1 thus leads to a belief-state expression of the input-output entropy. Specifically, the definition of $H_\mu(Y^T||U^{T-1})$ in (10) and the tower property of conditional expectation gives that

$$H_\mu(Y^T||U^{T-1}) = H(Y_0) + E_\mu \left[\sum_{k=0}^{T-1} \tilde{G}_3(\pi_k, U_k) \right] \quad (12)$$

where $\tilde{G}_3(\pi_k, u_k)$ is the entropy of the conditional pmf

$$p(y_{k+1}|\pi_k, u_k) = \sum_{x, \bar{x} \in \mathcal{X}} B^{x, y_{k+1}}(u_k) A^{x, \bar{x}}(u_k) \pi_k(\bar{x}), \quad (13)$$

that is, $H(Y_{k+1}|y^k, u^k)$, defined as

$$\tilde{G}_3(\pi_k, u_k) \triangleq - \sum_{y \in \mathcal{Y}} p(y|\pi_k, u_k) \log p(y|\pi_k, u_k). \quad (14)$$

A similar belief-state form of $H_\mu(U^{T-1}||Y^{T-1})$ also holds but will prove unnecessary since we shall next show that deterministic policies solve (7) (for which $H_\mu(U^{T-1}||Y^{T-1}) = 0$).

C. Belief MDP Reformulation

Along the lines of considering deterministic policies and omitting $H_\mu(U^{T-1}||Y^{T-1})$, the following lemma introduces a useful surrogate problem and is the final intermediate result we require to reformulate (7) as a belief MDP.

Lemma 3.2: If a deterministic policy $\mu^* \in \mathcal{P}$ minimizes

$$E_T[J_\mu^T + \beta H_\mu(X^T|Y^T, U^{T-1}) + \lambda H_\mu(Y^T||U^{T-1})] \quad (15)$$

over all policies $\mu \in \mathcal{P}$ under the same constraints as (7) given $\beta, \lambda \geq 0$, then μ^* also solves (7) with the same $\beta, \lambda \geq 0$.

Proof: The definition of the infimum implies that

$$\begin{aligned} & E_T[J_{\mu^*}^T + \beta H_{\mu^*}(X^T|Y^T, U^{T-1}) + \lambda H_{\mu^*}(Y^T||U^{T-1})] \\ & \geq \inf_{\mu \in \mathcal{P}} E_T[J_\mu^T + \beta H_\mu(X^T|Y^T, U^{T-1}) + \lambda H_\mu(Y^T||U^{T-1})] \\ & \geq \inf_{\mu \in \mathcal{P}} E_T[J_\mu^T + \beta H_\mu(X^T|Y^T, U^{T-1}) + \lambda H_\mu(Y^T||U^{T-1})] \\ & = E_T[J_{\mu^*}^T + \beta H_{\mu^*}(X^T|Y^T, U^{T-1}) + \lambda H_{\mu^*}(Y^T||U^{T-1})] \end{aligned}$$

where the second inequality holds due to Lemma 3.1 by noting that $H_\mu(U^{T-1}||Y^{T-1}) \geq 0$ for all $\mu \in \mathcal{P}$, and the last line holds via the definition of μ^* . These inequalities must hold with equality since Lemma 3.1 combined with $H_{\mu^*}(U^{T-1}||Y^{T-1}) = 0$ due to μ^* being deterministic implies

$$\begin{aligned} & E_T[J_{\mu^*}^T + \beta H_{\mu^*}(X^T|Y^T, U^{T-1}) + \lambda H_{\mu^*}(Y^T||U^{T-1})] \\ & = E_T[J_{\mu^*}^T + \beta H_{\mu^*}(X^T|Y^T, U^{T-1}) + \lambda H_{\mu^*}(Y^T||U^{T-1})]. \end{aligned}$$

The proof is complete. ■

A reformulation of (7) as a belief MDP follows.

Theorem 3.1: Define the belief-state cost function

$$\begin{aligned} G(\pi_k, u_k) & \\ & \triangleq (1 - \gamma)\beta\tilde{G}_1(\pi_k) + \gamma\beta\tilde{G}_2(\pi_k, u_k) + \gamma\lambda\tilde{G}_3(\pi_k, u_k) \\ & + E_{X_k}[(1 - \gamma)c_T(X_k) + \gamma c(X_k, U_k)|\pi_k, U_k = u_k]. \end{aligned} \quad (16)$$

Then (7) with $\beta, \lambda \geq 0$ is equivalent (up to $\lambda H(Y_0)$) to:

$$\begin{aligned} \inf_{\bar{\mu}} E_{\bar{\mu}} \left[\sum_{k=0}^{\infty} \gamma^k G(\pi_k, U_k) \middle| \pi_0 \right] \\ \text{s.t. } \pi_{k+1} = \Pi(\pi_k, U_k, Y_{k+1}) \\ Y_{k+1}|\pi_k, U_k \sim p(y_{k+1}|\pi_k, u_k) \\ U_k = \bar{\mu}(\pi_k) \in \mathcal{U} \end{aligned} \quad (17)$$

where the optimization is over deterministic, stationary policies $\bar{\mu} : \Delta \mapsto \mathcal{U}$ that are functions of the belief state π_k , and γ is the parameter of the geometric distribution of T .

Proof: Given Lemma 3.2, it suffices to show that minimizing (15) under the same constraints as (7) is equivalent (up to the constant $\lambda H(Y_0)$) to the belief MDP (17).

Rewriting (15) for any $\mu \in \mathcal{P}$ using (6), (8), and (12) gives

$$\begin{aligned} E_T[J_{\mu}^T + \beta H_{\mu}(X^T|Y^T, U^{T-1}) + \lambda H_{\mu}(Y^T||U^{T-1})] \\ = \lambda H(Y_0) + E_{T,\mu} \left[\tilde{G}_T(\pi_T) + \sum_{k=0}^{T-1} \tilde{G}(\pi_k, U_k) \right] \end{aligned}$$

via nested expectations with $\tilde{G}_T(\pi_T) \triangleq E_{X_T}[c_T(X_T) + \beta\tilde{G}_1(\pi_T)|\pi_T]$ and $\tilde{G}(\pi_k, U_k) \triangleq E_{X_k}[c(X_k, U_k) + \beta\tilde{G}_2(\pi_k, U_k) + \lambda\tilde{G}_3(\pi_k, U_k)|\pi_k, U_k]$. Ignoring $\lambda H(Y_0)$,

$$\begin{aligned} E_{T,\mu} \left[\tilde{G}_T(\pi_T) + \sum_{k=0}^{T-1} \tilde{G}(\pi_k, U_k) \right] \\ = E_{\mu} \left[\sum_{t=0}^{\infty} \zeta_t \left(\tilde{G}_T(\pi_t) + \sum_{k=0}^{t-1} \tilde{G}(\pi_k, U_k) \right) \right] \\ = E_{\mu} \left[\sum_{k=0}^{\infty} \left(\zeta_k \tilde{G}_T(\pi_k) + \sum_{t=k+1}^{\infty} \zeta_t \tilde{G}(\pi_k, U_k) \right) \right] \\ = E_{\mu} \left[\sum_{k=0}^{\infty} (\zeta_k \tilde{G}_T(\pi_k) + P(T > k) \tilde{G}(\pi_k, U_k)) \right] \\ = E_{\mu} \left[\sum_{k=0}^{\infty} \gamma^k \left((1 - \gamma) \tilde{G}_T(\pi_k) + \gamma \tilde{G}(\pi_k, U_k) \right) \right] \\ = E_{\mu} \left[\sum_{k=0}^{\infty} \gamma^k G(\pi_k, U_k) \right] \end{aligned}$$

where the second equality holds by interchanging summations; the third and fourth equalities follow from the cumulative distribution and pmf of the geometric distribution; and, the last equality holds by definition. Standard POMDP (or MDP) results imply that this expectation can be minimized over $\mu \in \mathcal{P}$ under the same constraints as (7) by deterministic stationary policies $\bar{\mu}$ that are functions of π_k (cf. [27, Section 5.4.1] and [10, Theorem 6.2.2]). The proof is complete. ■

D. Structural Results and Bounded-Error Solutions

Given the reformulation of (7) in Theorem 3.1, standard MDP or POMDP results (e.g., [10, Theorem 6.2.2] or [6]) imply that an optimal policy $\bar{\mu}^* : \Delta \mapsto \mathcal{U}$ and value function $V : \Delta \mapsto \mathbb{R}$ solving (7) satisfy Bellman's equation

$$V(\pi) = \min_{u \in \mathcal{U}} \{G(\pi, u) + \gamma E_Y [V(\Pi(\pi, u, Y)) | \pi, u]\} \quad (18)$$

for all $\pi \in \Delta$ with $\bar{\mu}^*(\pi)$ being a minimizing argument of (18). Solving (18) is, in general, difficult. However, if the functions G and V are concave in π , then standard POMDP techniques can yield solutions to (18). We thus examine G and V .

Theorem 3.2: The cost and value functions $G(\pi_k, u_k)$ and $V(\pi_k)$ of (7) reformulated as (17) are concave and continuous in $\pi_k \in \Delta$ for all $u_k \in \mathcal{U}$, all $0 < \gamma < 1$, and all $\beta, \lambda \geq 0$.

Proof: To prove the theorem assertion for G , it suffices to show that each function in (16) is concave and continuous in π_k since their coefficients are nonnegative for $\beta, \lambda \geq 0$ and $0 < \gamma < 1$. Firstly, \tilde{G}_1 and \tilde{G}_2 are concave and continuous in π_k via [5, Lemma 2]. Secondly, \tilde{G}_3 is the entropy of the conditional pmf $p(y|\pi_k, u_k)$, and so is concave and continuous in it via [23, Theorem 2.7.3]. Since $p(y|\pi_k, u_k)$ is linear in π_k (cf. (13)), \tilde{G}_3 is concave in a linear function of π_k , and so is concave and continuous in π_k . Finally, the expectation in (16) is concave and continuous in π_k since it equals $\sum_{x \in \mathcal{X}} \pi_k(x) [(1 - \gamma)c_T(x) + \gamma c(x, u_k)]$. That V is concave and continuous follows via [6, Theorem 3.1]. ■

Theorem 3.2 enables the use of standard POMDP techniques to find bounded-error solutions to (7). Specifically, following the PWLC approach proposed in [6], consider an arbitrary finite set $\Xi \subset \Delta$ of *base points* $\xi \in \Xi$ at which the gradient $\nabla_\pi G(\xi, u)$ of $G(\cdot, u)$ is well defined for all $u \in \mathcal{U}$. For each $u \in \mathcal{U}$, the tangent hyperplane to $G(\cdot, u)$ at each $\xi \in \Xi$ is

$$\omega_\xi^u(\pi) \triangleq G(\xi, u) + \langle (\pi - \xi), \nabla_\pi G(\xi, u) \rangle = \langle \pi, \alpha_\xi^u \rangle$$

for $\pi \in \Delta$ where $\langle \cdot, \cdot \rangle$ denotes the inner product, and $\alpha_\xi^u \triangleq G(\xi, u) + \nabla_\pi G(\xi, u) - \langle \xi, \nabla_\pi G(\xi, u) \rangle \in \mathbb{R}^{N_x}$ are vectors (with the addition of a vector and a scalar here meaning the addition of the scalar

to all components of the vector). Since G is concave via Theorem 3.2, the hyperplanes form a PWLC approximation \hat{G} to G ; that is, for $\pi \in \Delta$ and $u \in \mathcal{U}$,

$$\hat{G}(\pi, u) \triangleq \min_{\xi \in \Xi} \langle \pi, \alpha_\xi^u \rangle \geq G(\pi, u).$$

By replacing G in (18) with \hat{G} , (18) can be solved for an approximate PWLC value function \hat{V} (and policy) using standard POMDP algorithms that operate directly on the vectors $\{\alpha_\xi^u : \xi \in \Xi, u \in \mathcal{U}\}$ (see [6, Section 3.3] for more details). Furthermore, G satisfies the Hölder continuity condition of [6, Theorem 4.3] since the (negative) entropy function $f(x) = \sum_{i=1}^{N_x} x(i) \log x(i)$ is Hölder continuous on Δ (as are continuous linear functions, and the sums and compositions of Hölder continuous functions, cf. [28, Example 1.1.4, and Propositions 1.2.1 and 1.2.2]). Hence, [6, Section 4.2] implies that there exists constants $\kappa_1 > 0$ and $\kappa_2 \in (0, 1)$ such that $\|V - \hat{V}\|_\infty \leq \kappa_1 (\delta_\Xi)^{\kappa_2}$ where $\delta_\Xi \triangleq \min_{\pi \in \Delta} \max_{\xi \in \Xi} \|\pi - \xi\|_1$ is the sparsity of Ξ with $\|\cdot\|_1$ and $\|\cdot\|_\infty$ denoting the l^1 -norm and L^∞ -norm, respectively. In principle, this error can be made arbitrarily small by selecting $\xi \in \Xi$ to decrease δ_Ξ .

IV. SPECIAL CASE OF JOINT-ENTROPY REGULARIZATION

The PWLC approach to solving (7) presented in the previous section is consistent with state-of-the-art approaches to solving POMDPs with nonlinear belief-state cost functions (see [10, Chapter 8], [6], [7]). However, constructing accurate PWLC approximations can require a large number of linear segments (i.e. vectors α_ξ^u), resulting in significant computational effort and the need to modify standard POMDP solver implementations (cf. [10, Section 8.4.5]). In this section, we explore a simpler approach to solving (7) that is tractable *without* PWLC approximations when the smoother and input-output entropies are equally penalized, that is, when $\beta = \lambda \geq 0$ so that the sum of the smoother and input-output entropies in (7) becomes the *joint entropy* defined as $H_\mu(X^T, Y^T, U^{T-1}) \triangleq -E_\mu[\log p_\mu(X^T, Y^T, U^{T-1})] = H_\mu(X^T | Y^T, U^{T-1}) + H_\mu(Y^T, U^{T-1})$. Key to this simpler approach is the following new expression for the joint entropy.

Lemma 4.1: For any policy $\mu \in \mathcal{P}$ and $T \geq 0$, we have:

$$\begin{aligned} & H_\mu(X^T, Y^T, U^{T-1}) \\ &= H(X_0, Y_0) + H_\mu(U^{T-1} \| Y^{T-1}) + E_\mu \left[\sum_{k=0}^{T-1} \tilde{c}(X_k, U_k) \right] \end{aligned}$$

where

$$\begin{aligned} & \tilde{c}(x_k, u_k) \\ & \triangleq - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} A^{x, x_k}(u_k) B^{x, y}(u_k) \log(A^{x, x_k}(u_k) B^{x, y}(u_k)). \end{aligned}$$

Proof: From the definition of the joint entropy and (3),

$$\begin{aligned}
& H_\mu(X^T, Y^T, U^{T-1}) \\
&= -E_\mu \left[\log(\rho(X_0)B^{X_0, Y_0}) + \sum_{k=0}^{T-1} \log(\mu_k^{I_k}(U_k)Z_k) \right] \\
&= H(X_0, Y_0) - \sum_{k=0}^{T-1} E_\mu [\log \mu_k^{I_k}(U_k)] - \sum_{k=0}^{T-1} E_\mu [\log Z_k] \\
&= H(X_0, Y_0) + H_\mu(U^{T-1} \| Y^{T-1}) + E_\mu \left[\sum_{k=0}^{T-1} \tilde{c}(X_k, U_k) \right]
\end{aligned}$$

where $Z_k \triangleq A^{X_{k+1}, X_k}(U_k)B^{X_{k+1}, Y_{k+1}}(U_k)$; the second equality holds due to the properties of the logarithm and linearity of expectations and summations; and, the third equality follows from (11), with nested expectations giving

$$-E_\mu [\log Z_k] = E_\mu [E_{X_{k+1}, Y_{k+1}} [-\log Z_k | X_k, U_k]]$$

where the inner expectation is \tilde{c} . The proof is complete. \blacksquare

A second reformulation of (7) with $\beta = \lambda \geq 0$ follows.

Theorem 4.1: Define the belief-state cost function

$$L(\pi_k, u_k) \triangleq E_{X_k}[\ell(X_k, u_k) | \pi_k, u_k] = \sum_{x \in \mathcal{X}} \pi_k(x) \ell(x, u_k)$$

where $\ell(x_k, u_k) \triangleq (1 - \gamma)c_T(x_k) + \gamma c(x_k, u_k) + \gamma\beta\tilde{c}(x_k, u_k)$, then (7) with $\beta = \lambda \geq 0$ is equivalent (up to $\beta H(X_0, Y_0)$) to:

$$\inf_{\bar{\mu}} E_{\bar{\mu}} \left[\sum_{k=0}^{\infty} \gamma^k L(\pi_k, U_k) \middle| \pi_0 \right] \quad (19)$$

subject to the same constraints as (17) and where the optimization is over deterministic, stationary policies $\bar{\mu} : \Delta \mapsto \mathcal{U}$.

Proof: Same as that of Theorem 3.1, but using Lemma 4.1 instead of (8) and (12) to rewrite (15), noting that $H_\mu(X^T | Y^T, U^{T-1}) + H_\mu(Y^T \| U^{T-1}) = H_\mu(X^T, Y^T, U^{T-1}) - H_\mu(U^{T-1} \| Y^{T-1})$ via Lemma 3.1. \blacksquare

The belief MDP reformulation of our ERPOMDP problem in (19) with $\beta = \lambda$ is surprising because its cost function L is *linear* in π_k . In contrast, the cost function G of the first belief MDP reformulation established in (17) is *nonlinear* in π_k , even when $\beta = \lambda$. The two different belief MDP reformulations of (7) in (17) and (19) are due to the joint pmf $p_\mu(x^T, y^T, u^{T-1})$ admitting multiple factorizations, with (3) leading to (19), and a factorization similar to (9) leading to (17).

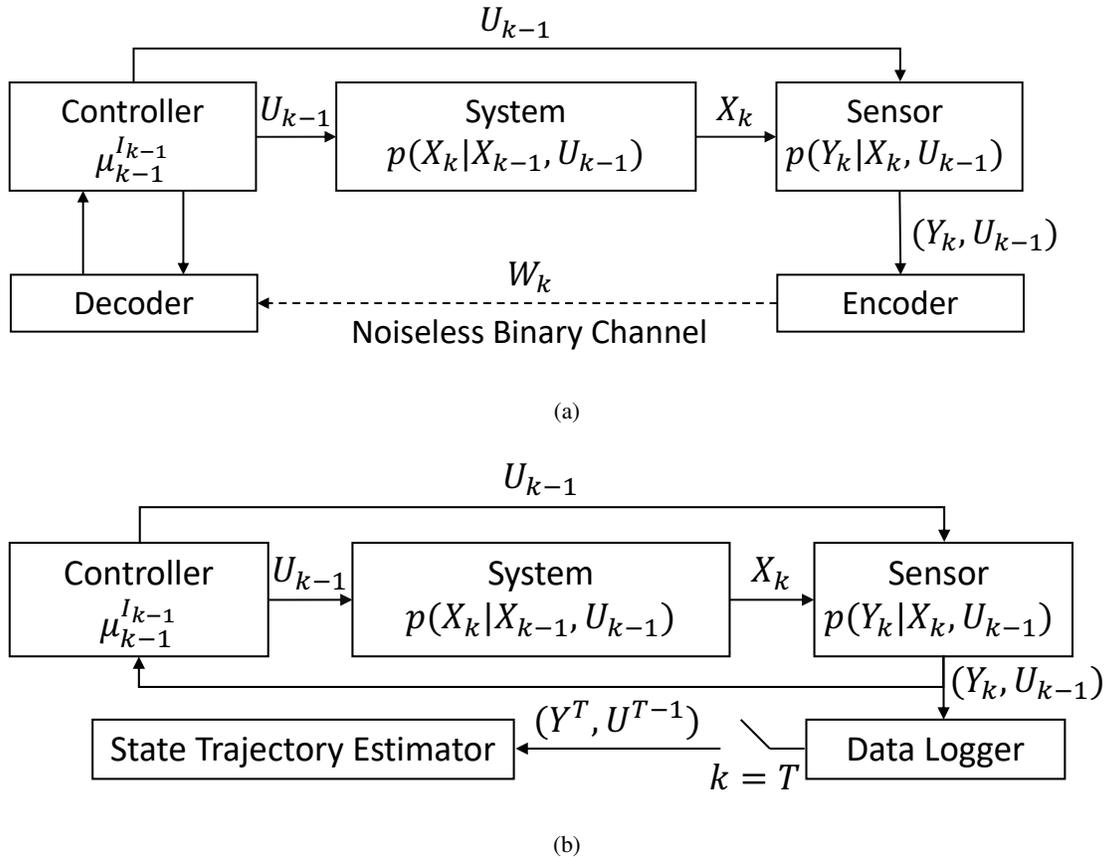


Fig. 1. Operational Interpretations of the ERPOMDP problem (7): (a) Networked Control; and (b) Memory-Efficient Active State Trajectory Estimation.

The linearity of L is of considerable practical value because it enables (7) with $\beta = \lambda$ to be solved using standard POMDP solution techniques without any PWLC approximation of L . Indeed, for any $u \in \mathcal{U}$ and $\pi \in \Delta$, $L(\pi, u) = \langle \pi, \alpha^u \rangle$ holds exactly given the (single) vector $\alpha^u \triangleq [\ell(1, u), \ell(2, u), \dots, \ell(N, u)]'$. Dynamic programming equations of the form of (18) with L replacing G can thus be solved using standard POMDP techniques that operate directly on the vectors $\{\alpha^u : u \in \mathcal{U}\}$ (cf. [10, Chapter 7.5]). We next discuss the operational significance of the linearity of L .

V. OPERATIONAL INTERPRETATIONS AND RELATIONSHIPS

In this section, we discuss two operational interpretations of ERPOMDPs, and discuss their relationship to other optimization problems with information-theoretic terms.

A. Networked Control Rate-Cost Trade-Offs

Consider a networked control setting in which the feedback path of a POMDP involves transmission over a noiseless binary channel, as illustrated in Fig. 1a. At every time step k , an encoder receives observations Y_k of the state X_k (e.g., arising from a sensor), and selects and transmits a binary codeword W_k from a predefined codebook of optimal uniquely decodable binary codes \mathcal{W}_k . Upon receiving W_k , the decoder decodes Y_k and passes it to a controller, that uses it to evaluate the next control U_k . We allow the codebook \mathcal{W}_k to be time-varying, and let R_k be the length (number of bits) of the codeword W_k . We assume that both the encoder and decoder have infinite memories so that W_k conveys Y_k given (Y^{k-1}, U^{k-1}) . Shannon's source coding theorem (cf. [23, Section 5.5]) then implies that the (minimum) expected data transmitted in the noiseless binary channel at time k satisfies

$$H_\mu(Y_k|Y^{k-1}, U^{k-1}) \leq E_\mu[R_k] \leq H_\mu(Y_k|Y^{k-1}, U^{k-1}) + 1,$$

and so the total expected data transmitted over T satisfies

$$\begin{aligned} E_T[H_\mu(Y^T||U^{T-1})] &\leq E_{T,\mu} \left[\sum_{k=0}^T R_k \right] \\ &\leq E_T [H_\mu(Y^T||U^{T-1}) + T + 1]. \end{aligned}$$

In view of Lemma 3.2 and Theorem 3.1, our ERPOMDP problem (7) with $\lambda > 0$ (and any $\beta \geq 0$) thus has the operational interpretation of seeking policies that trade-off the expected total data transmitted for feedback control via $H_\mu(Y^T||U^{T-1})$, with the value of the cost functional J_μ^T .

Theorem 3.1 is particularly important for introducing rate-cost trade-offs into POMDPs because it enables regularization only by the input-output entropy (i.e., it holds when $\lambda > 0$ but $\beta = 0$ in (7)). In contrast, Theorem 4.1 is, in general, of secondary value for introducing rate-cost trade-offs because it also requires regularization by the smoother entropy (i.e., it holds only when $\beta = \lambda$ in (7)), but does enable the use of standard POMDP techniques without PWLC approximations.

B. Memory-Efficient Active State Trajectory Estimation

The second operational interpretation of our ERPOMDP problem (7) relates to active state estimation (i.e. controlling a POMDP to aid the estimation of its latent state trajectory X^T from stored trajectories (Y^T, U^{T-1})). Such problems arise in robotics (cf. [13], [14]) and controlled sensing (cf. [9], [10]).

As shown in Fig. 1b, consider a POMDP in which, at each time k , the current observation and control (Y_k, U_{k-1}) are encoded and stored by a data logger by selecting and storing a binary codeword W_k from a predefined codebook of optimal uniquely decodable binary codes \mathcal{W}_k . Thus, the data logger encodes

(Y_k, U_{k-1}) given (Y^{k-1}, U^{k-2}) . We allow the codebook \mathcal{W}_k to be time-varying, and let R_k be the length of the codeword W_k . Shannon's source coding theorem (cf. [23, Section 5.5]) implies that the (minimum) expected total memory required to store (Y^T, U^{T-1}) satisfies the bounds

$$\begin{aligned} E_T[H_\mu(Y^T, U^{T-1})] &\leq E_{T,\mu} \left[\sum_{k=0}^T R_k \right] \\ &\leq E_T [H_\mu(Y^T, U^{T-1}) + T + 1]. \end{aligned}$$

At the conclusion of the control horizon ($k = T$), the data logger decodes the stored trajectories (Y^T, U^{T-1}) and passes them to an (offline) algorithm for estimating the state trajectory X^T (e.g. the Viterbi algorithm). Let the minimum probability of error for *any* estimator of X^T given (Y^T, U^{T-1}) be $\epsilon \triangleq \min_{\hat{X}^T} P(X^T \neq \hat{X}^T)$ with \hat{X}^T being any function $f : \mathcal{Y}^T \times \mathcal{U}^{T-1} \mapsto \mathcal{X}^T$. Theorem 1 of [29] gives that

$$\Phi^{-1}(H_\mu(X^T|Y^T, U^{T-1})) \leq \epsilon \leq \phi^{-1}(H_\mu(X^T|Y^T, U^{T-1}))$$

where Φ^{-1} and ϕ^{-1} are the inverse functions of strictly monotonically increasing functions (defined in [29]), and thus are themselves strictly monotonically increasing.

The involvement of the smoother and input-output entropies in bounds on the estimation error and required memory implies that (7) has the operational interpretation of seeking policies that aid estimation of the state trajectory (by reducing the smoother entropy) whilst decreasing the memory required to store the observation and control trajectories (by reducing the input-output entropy). In the special case considered in Section IV with $\beta = \lambda$, (7) constitutes a formulation of active state estimation in which the estimation and memory objectives are weighted equally. In this regard, Theorem 4.1 establishing the linearity of L in (19) is further surprising since most previous active state estimation formulations involve cost that are entirely nonlinear in the belief state and can only be optimized by resorting to approximations (cf. [6], [9], [10]).

C. Relationship to Other Information-Theoretic POMDPs

ERPOMDPs (7) are closely related to problems involving the optimization of information-theoretic terms that have previously been considered for reinforcement learning [30], studying the capacity of channels with memory and feedback [31], privacy (e.g. in smart metering systems) [19], [21], [32], [33], and studying rate-cost trade-offs in MDPs and POMDPs [2], [8]. These problems, however, mostly involve optimizing only a single information-theoretic term derived from either the mutual information between states and/or observations (e.g., directed information and transfer entropy [2], [8], [19], [21], [31], [32]), or the entropy of the states or controls (e.g., $H_\mu(U_k|i_k) = -\sum_{u \in \mathcal{U}} \mu_k^{i_k}(u) \log \mu_k^{i_k}(u)$ [30], [33]). In

contrast, ERPOMDPs involve both the standard cost functional J_μ^T and, in general, two information-theoretic terms, the smoother entropy and the (novel) input-output entropy. The procedure for solving ERPOMDPs is, however, similar to that of solving these other optimization problems, with most having been shown to lead to belief MDPs — albeit few (if any) with cost functions that are linear in the belief state, rendering our Theorem 4.1 joint-entropy result further surprising.

VI. SIMULATION EXAMPLE

We now simulate ERPOMDPs for active state estimation.

A. Example Set-Up

Consider an agent in the grid shown in Fig. 2, that seeks to move to (and stay in) a known goal location from an unknown starting location (distributed uniformly over the grid such that the initial state pmf ρ is uniform), whilst actively localizing itself so as to enable its path to the goal to be estimated for the purpose of later being retraced or communicated. Each cell in the grid is a state in the agent’s state space $\mathcal{X} = \{1, \dots, 144\}$ (enumerated top-to-bottom, left-to-right). The agent has five possible control actions $\mathcal{U} = \{1, \dots, 5\}$, corresponding to moving one cell in each of the four compass directions, or staying still (all with probability 1). There are internal and external walls (bold black lines in Fig. 2) that block movement, with the agent staying still if it attempts to move into them. The agent receives measurements $\mathcal{Y} = \{1, \dots, 16\}$ corresponding to whether or not a wall is immediately adjacent to its current cell in each of the four compass directions. The agent detects a wall when it is present (resp. not present) with probability 1 (resp. 0.2). A simplified version of this example was previously considered in [2] for MDPs.

We examine the ability of the agent to move to the goal and ensure estimation of its path by solving (7) with either $\beta = \lambda = 0$ (corresponding to a standard POMDP without any regularization), $\beta = \lambda = 1$ (corresponding to joint-entropy regularization), $\beta = 1$ and $\lambda = 0$ (corresponding to only smoother-entropy regularization), and $\beta = 0$ and $\lambda = 1$ (corresponding to only input-output-entropy regularization). In all cases, $\gamma = 0.99$ and the goal objective is encoded via the cost $c(x, u) = 1_{\{x \neq 144\}}$ for all $x \in \mathcal{X}$ and $u \in \mathcal{U}$.

We use SARSOP [24] to solve (7) via the reformulation in Theorem 4.1 when $\beta = \lambda \in \{0, 1\}$, and via the reformulation in Theorem 3.1 when $\beta \neq \lambda$. For $\beta \neq \lambda$, we construct a PWLC approximation of G in (17) using a set Ξ containing the middle of the simplex Δ and points near the vertices with values in their largest element of 0.857 and 0.001 in their other 143 elements. For $\beta = \lambda$, we avoid PWLC approximations since the cost function L in Theorem 4.1 is linear. From Table I, we see that the time

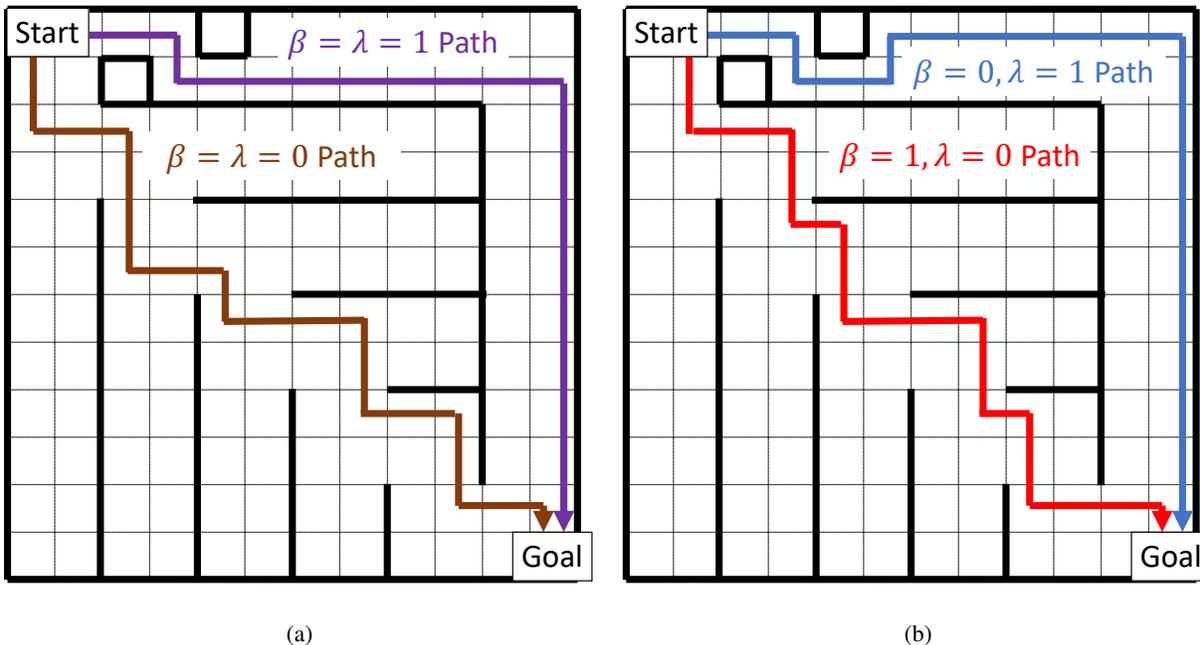


Fig. 2. Example realizations of ERPOMDP (7) agent with wall sensor moving from start cell (unknown to agent, top-left in these realizations) to goal (walls bold black): (a) $\beta = \lambda \in \{0, 1\}$ (b) $\beta = 1$ & $\lambda = 0$, and $\beta = 0$ & $\lambda = 1$.

taken to compute policies requiring PWLC approximations (i.e., policies with $\beta \neq \lambda$) is much greater than the time required to compute the standard POMDP policy with $\beta = \lambda = 0$, and the ERPOMDP policy with $\beta = \lambda = 1$.

B. Simulation Results

The results of 1000 Monte Carlo simulations of each policy over $T = 100$ time steps (the mean of T) are summarized in Table I. For each policy, we report several active state estimation criteria including: the total (undiscounted) cost associated with not being in the goal state (i.e., the *Goal Cost*) $\sum_{k=0}^T E[c(X_k, U_k)]$; the input-output entropy; the smoother entropy; the joint entropy; the sum of belief entropies $\sum_{k=0}^T H_\mu(X_k | Y^k, U^{k-1})$; and, the probability of error in maximum *a posteriori* estimates of the trajectory X^T (*Traj. MAP Error Prob.*) computed via the Viterbi algorithm [10, Section 3.5.3]. Example state realizations with the agent starting in the top-left cell are shown in Fig. 2.

Table I shows that the standard POMDP policy ($\beta = \lambda = 0$) results in the lowest goal cost, which is unsurprising since it only explicitly minimizes the (discounted) cost of the agent not being in the goal state. The smoother entropy, input-output entropy, and joint entropy are all significantly less (better) when they are regularized via selection of $\beta = 1$ and $\lambda = 0$, $\lambda = 1$ and $\beta = 0$, or both $\beta = \lambda = 1$, respectively (at the expense of a small increase in the *Goal Cost*). Table I thus highlights that the agent resolves its

TABLE I
 MONTE CARLO SIMULATION RESULTS (BEST VALUES IN BOLD). COMPUTATIONAL TIMES FOR AN M1 2020 APPLE
 MACBOOK AIR.

Performance Criteria	ERPOMDP (7) Policy Parameters			
	$\beta = 0$	$\beta = 1$	$\beta = 0$	$\beta = 1$
	$\lambda = 0$	$\lambda = 0$	$\lambda = 1$	$\lambda = 1$
Goal Cost	23.0	25.2	25.9	26.0
Input-Output Entropy	120.2	122.0	114.9	114.9
Smoother Entropy	1.47	0.41	0.60	0.50
Joint Entropy	121.7	122.4	115.5	115.4
Sum of Belief Entropies	21.9	17.3	16.0	16.0
Traj. MAP Error Prob.	0.01	0.00	0.01	0.00
Time to Compute Policy (s)	0.21	6563	964	0.66

uncertainty and reduces the memory required to store its measurements more effectively with versions of ERPOMDP policies with nonzero β and λ than with the standard POMDP policy (with $\beta = \lambda = 0$).

As illustrated in Fig. 2, the ERPOMDP policies with $\beta = \lambda = 1$, and $\beta = 0$ and $\lambda = 1$, reduce the joint entropy and input-output entropy most effectively by moving the agent through the cells along the walls since these yield easily compressible observation sequences (a single wall in the same relative direction). In contrast, the cells in the center through which the standard POMDP policy (with $\beta = \lambda = 0$) and the ERPOMDP policy with $\beta = 1$ and $\lambda = 0$ move the agent yield more variable (higher entropy) observations since they are surrounded by 0, 1, or 2 walls. By initially moving the agent East from its starting cell and then through the center, the ERPOMDP policy with $\beta = 1$ and $\lambda = 0$ is able to achieve the minimum smoother entropy. However, in this example, the differences in smoother entropies between policies with $\beta = 0$ and $\beta = 1$ is much smaller compared to the differences in input-output entropies between policies with $\lambda = 0$ and $\lambda = 1$. Thus, the ERPOMDP policy involving joint-entropy regularization via $\beta = \lambda = 1$ achieves close to the best performance across the criteria (at the slight expense of the *Goal Cost*) whilst avoiding PWLC approximations, highlighting the value of Theorem 4.1. Clearly, finer trade-offs between the *Goal Cost* and smoother or input-output entropies can be obtained by selecting β and λ independently and using Theorem 3.1, but at considerable computational cost.

VII. CONCLUSION

We propose ERPOMDPs, show that they admit PWLC approximate solutions, and discuss their relevance to active state estimation and rate-cost trade-offs. Surprisingly, ERPOMDPs admit exact solutions

when regularizing by the joint entropy of the states, observations, and controls, which constitutes a novel, tractable formulation of active state estimation.

REFERENCES

- [1] T. L. Molloy and G. N. Nair, “JEM: Joint Entropy Minimization for Active State Estimation with Linear POMDP Costs,” in *2022 American Control Conference (ACC)*, 2022, pp. 1601–1607.
- [2] T. Tanaka, H. Sandberg, and M. Skoglund, “Transfer-Entropy-Regularized Markov Decision Processes,” *IEEE Transactions on Automatic Control*, pp. 1–1, 2021.
- [3] E. Shafieepoorfard, M. Raginsky, and S. P. Meyn, “Rationally inattentive control of Markov processes,” *SIAM Journal on Control and Optimization*, vol. 54, no. 2, pp. 987–1016, 2016.
- [4] G. M. Hoffmann and C. J. Tomlin, “Mobile sensor network control using mutual information methods and particle filters,” *IEEE Transactions on Automatic Control*, vol. 55, no. 1, pp. 32–47, 2010.
- [5] T. L. Molloy and G. N. Nair, “Active trajectory estimation for partially observed Markov decision processes via conditional entropy,” in *2021 European Control Conference (ECC)*, 2021, pp. 385–391.
- [6] M. Araya, O. Buffet, V. Thomas, and F. Charpillat, “A POMDP extension with belief-dependent rewards,” in *Advances in Neural Information Processing Systems*, J. Lafferty, C. Williams, J. Shawe-Taylor, R. Zemel, and A. Culotta, Eds., vol. 23. Curran Associates, Inc., 2010, pp. 64–72.
- [7] M. Fehr, O. Buffet, V. Thomas, and J. Dibangoye, “rho-POMDPs have Lipschitz-Continuous epsilon-Optimal Value Functions,” in *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., vol. 31. Curran Associates, Inc., 2018.
- [8] J. Rubin, O. Shamir, and N. Tishby, “Trading value and information in MDPs,” in *Decision Making with Imperfect Decision Makers*. Springer, 2012, pp. 57–74.
- [9] V. Krishnamurthy and D. V. Djonin, “Structured threshold policies for dynamic sensor scheduling – a partially observed Markov decision process approach,” *IEEE Transactions on Signal Processing*, vol. 55, no. 10, pp. 4938–4957, 2007.
- [10] V. Krishnamurthy, *Partially observed Markov decision processes*. Cambridge University Press, 2016.
- [11] D.-S. Zois and U. Mitra, “Active state tracking with sensing costs: Analysis of two-states and methods for n -states,” *IEEE Transactions on Signal Processing*, vol. 65, no. 11, pp. 2828–2843, 2017.
- [12] D.-S. Zois, M. Levorato, and U. Mitra, “Active classification for POMDPs: A Kalman-like state estimator,” *IEEE Transactions on Signal Processing*, vol. 62, no. 23, pp. 6209–6224, 2014.
- [13] S. Thrun, W. Burgard, and D. Fox, *Probabilistic Robotics*. MIT Press, 2005.
- [14] C. Stachniss, G. Grisetti, and W. Burgard, “Information gain-based exploration using Rao-Blackwellized particle filters.” in *Robotics: Science and Systems*, vol. 2, 2005, pp. 65–72.
- [15] R. Valencia, J. Valls Miró, G. Dissanayake, and J. Andrade-Cetto, “Active Pose SLAM,” in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2012, pp. 1885–1891.
- [16] T. L. Molloy and G. N. Nair, “Smoother Entropy for Active State Trajectory Estimation and Obfuscation in POMDPs,” *arXiv preprint arXiv:2108.10227*, 2021.
- [17] E. G. Ryan, C. C. Drovandi, J. M. McGree, and A. N. Pettitt, “A review of modern computational algorithms for Bayesian optimal design,” *International Statistical Review*, vol. 84, no. 1, pp. 128–154, 2016.
- [18] K. Chaloner and I. Verdinelli, “Bayesian experimental design: A review,” *Statistical Science*, pp. 273–304, 1995.
- [19] T. Tanaka, P. M. Esfahani, and S. K. Mitter, “LQG control with minimum directed information: Semidefinite programming approach,” *IEEE Trans. on Automatic Control*, vol. 63, no. 1, pp. 37–52, 2018.

- [20] V. Kostina and B. Hassibi, “Rate-cost tradeoffs in control,” *IEEE Trans. on Automatic Control*, vol. 64, no. 11, pp. 4525–4540, 2019.
- [21] O. Sabag, P. Tian, V. Kostina, and B. Hassibi, “The Minimal Directed Information Needed to Improve the LQG Cost,” in *2020 59th IEEE Conference on Decision and Control (CDC)*, 2020, pp. 1842–1847.
- [22] A. Shwartz, “Death and discounting,” *IEEE Transactions on Automatic Control*, vol. 46, no. 4, pp. 644–647, 2001.
- [23] T. Cover and J. Thomas, *Elements of information theory*, 2nd ed. New York: Wiley, 2006.
- [24] H. Kurniawati, D. Hsu, and W. S. Lee, “SARSOP: Efficient point-based POMDP planning by approximating optimally reachable belief spaces.” in *Robotics: Science and Systems*, 2008.
- [25] N. P. Garg, D. Hsu, and W. S. Lee, “DESPOT-Alpha: Online POMDP planning with large state and observation spaces.” in *Robotics: Science and Systems*, 2019.
- [26] G. Kramer, “Directed information for channels with feedback,” Ph.D. dissertation, Swiss Federal Institute of Technology, Zurich, 1998.
- [27] D. P. Bertsekas, *Dynamic programming and optimal control*, Third ed. Belmont, MA: Athena Scientific, 1995, vol. 1.
- [28] R. Fiorenza, *Hölder and locally Hölder Continuous Functions, and Open Sets of Class C^k , $C^{k,\lambda}$* . Birkhäuser, 2017.
- [29] M. Feder and N. Merhav, “Relations between entropy and error probability,” *IEEE Trans. on Info. Theory*, vol. 40, no. 1, pp. 259–266, 1994.
- [30] T. Haarnoja, H. Tang, P. Abbeel, and S. Levine, “Reinforcement learning with deep energy-based policies,” in *International Conference on Machine Learning*. PMLR, 2017, pp. 1352–1361.
- [31] C. D. Charalambous and P. A. Stavrou, “Directed information on abstract spaces: Properties and variational equalities,” *IEEE Transactions on Information Theory*, vol. 62, no. 11, pp. 6019–6052, 2016.
- [32] S. Li, A. Khisti, and A. Mahajan, “Information-theoretic privacy for smart metering systems with a rechargeable battery,” *IEEE Trans. on Information Theory*, vol. 64, no. 5, pp. 3679–3695, 2018.
- [33] Y. Savas, M. Ornik, M. Cubuktepe, M. O. Karabag, and U. Topcu, “Entropy maximization for Markov decision processes under temporal logic constraints,” *IEEE Trans. on Automatic Control*, vol. 65, no. 4, pp. 1552–1567, 2020.