

# Principled Learning-to-Communicate with Quasi-Classical Information Structures

Xiangyu Liu<sup>†</sup>

Haoyi You<sup>†</sup>

Kaiqing Zhang<sup>†</sup>

**Abstract**—Learning-to-Communicate (LTC) in partially observable environments has emerged and received increasing attention in deep multi-agent reinforcement learning, where the control and communication strategies are *jointly* learned. On the other hand, the impact of communication has been extensively studied in control theory. In this paper, we seek to formalize and better understand LTC by bridging these two lines of work, through the lens of *information structures* (ISs). To this end, we formalize LTC in decentralized partially observable Markov decision processes (Dec-POMDPs) under the common-information-based framework from decentralized stochastic control, and classify LTC problems based on the ISs before (additional) information sharing. We first show that non-classical LTCs are computationally intractable in general, and thus focus on quasi-classical (QC) LTCs. We then propose a series of conditions for QC LTCs, violating which can cause computational hardness in general. Further, we develop provable planning and learning algorithms for QC LTCs, and show that some examples of QC LTCs satisfying the above conditions can be solved with quasi-polynomial time and samples. Along the way, we also establish some relationship between (strictly) QC IS and the condition of having strategy-independent common-information-based beliefs (SI-CIBs), as well as solving Dec-POMDPs without computationally intractable oracles but beyond those with the SI-CIB condition, which may be of independent interest.

## I. INTRODUCTION

The Learning-to-Communicate (LTC) problem has emerged and gained traction in the area of (deep) multi-agent reinforcement learning (MARL) [2], [3], [4]. Unlike classical MARL, which aims to learn a *control* strategy that minimizes the expected accumulated costs, LTC seeks to *jointly* minimize over both the *control* and the *communication* strategies of all the agents, as a way to mitigate the challenges due to the agents' *partial observability* of the environment. Despite the promising empirical successes, theoretical understandings of LTC remain largely underexplored.

On the other hand, in control theory, a rich literature has investigated the role of *communication* in decentralized/networked control [5], [6], inspiring us to rigorously examine LTCs from such a principled perspective. Most of these studies, however, focused on linear systems, and did not explore the *non-asymptotic* computational and/or sample complexity guarantees when the system knowledge is not known. A few recent studies [7], [8] started to explore the

settings with general discrete (nonlinear) spaces, with special communication protocols and state transition dynamics.

More broadly, the design of communication strategies dictates the *information structure* (IS) of the control system, which characterizes *who knows what and when* [9]. IS and its impact on the *optimization tractability*, especially for linear systems, have been extensively studied in decentralized control, see [10], [11] for comprehensive overviews. In this work, we seek a more principled understanding of LTCs through the lens of information structures, with a focus on the computational and sample complexities of the problem.

Specifically, we formalize LTCs in the general framework of decentralized partially observable Markov decision processes (Dec-POMDPs) [12], as in the empirical work [2], [3], [4], and study its finite-time and sample complexity guarantees. We detail our contributions as follows.

**Contributions.** (i) We formalize Learning-to-Communicate in Dec-POMDPs under the common-information-based framework from decentralized stochastic control [13], [14], [15], allowing the sharing of *historical* information and communication costs; (ii) We classify LTCs through the lens of information structures, according to the ISs before (additional) information sharing. We then show that LTCs with *non-classical* [10] baseline IS can be computationally intractable; (iii) Given the hardness, we focus on *quasi-classical* (QC) LTCs, and propose a series of conditions under which LTCs preserve the QC IS after sharing, while violating which can cause computational hardness in general; (iv) We propose both planning and learning algorithms for QC LTCs, by reformulating them as Dec-POMDPs with *strategy-independent common-information-based beliefs* (SI-CIBs) [14], [15], a condition previously shown to be critical for taming computational intractability [15]; (v) Quasi-polynomial time and sample complexities of the algorithms are established for QC LTC examples that satisfy the conditions in (iii). Along the way, we also establish some relationship between (strictly) quasi-classical ((s)QC) ISs and the SI-CIB condition, as well as solving general Dec-POMDPs without computationally intractable oracles but beyond the SI-CIB ones, thus advancing the results in [15].

## II. PRELIMINARIES

### A. Learning-to-Communicate (with Communication Cost)

For  $n > 1$  agents, a (cooperative) *Learning-to-Communicate* problem is described by a tuple  $\mathcal{L} = \langle H, \mathcal{S}, \{\mathcal{A}_{i,h}\}_{i \in [n], h \in [H]}, \{\mathcal{O}_{i,h}\}_{i \in [n], h \in [H]}, \{\mathcal{M}_{i,h}\}_{i \in [n], h \in [H]}, \mathbb{T}, \mathbb{O}, \mu_1, \{\mathcal{R}_h\}_{h \in [H]}, \{\mathcal{K}_h\}_{h \in [H]} \rangle$ , where  $H$  denotes the length of each episode. Other components are specified as follows.

<sup>†</sup>The authors are ordered alphabetically, and are affiliated with the University of Maryland, College Park, MD, USA, 20742. Emails: {xyliu999, yuriiyou, kaiqing}@umd.edu. This work was supported by the Army Research Office (ARO) grant W911NF-24-1-0085 and the NSF CAREER Award 2443704. A comprehensive technical report that contains all the omitted details can be found at [1].

1) *Decision-making components*: We use  $\mathcal{S}$  to denote the state space, and  $\mathcal{A}_{i,h}$  to denote the *control action* space of agent  $i$  at timestep  $h \in [H]$ . We denote by  $s_h \in \mathcal{S}$  the state and by  $a_{i,h}$  the control action of agent  $i$  at timestep  $h$ . We use  $a_h := (a_{1,h}, \dots, a_{n,h}) \in \mathcal{A}_h := \prod_{i \in [n]} \mathcal{A}_{i,h}$  to denote the joint control action of all the agents at timestep  $h$ . We denote by  $\mathbb{T} = \{\mathbb{T}_h\}_{h \in [H]}$  the collection of state transition kernels, where  $s_{h+1} \sim \mathbb{T}_h(\cdot | s_h, a_h) \in \Delta(\mathcal{S})$ . We use  $\mu_1 \in \Delta(\mathcal{S})$  to denote the initial state distribution,  $o_{i,h} \in \mathcal{O}_{i,h}$  to denote the observation of agent  $i$  at timestep  $h$ , and  $o_h := (o_{1,h}, o_{2,h}, \dots, o_{n,h}) \in \mathcal{O}_h := \mathcal{O}_{1,h} \times \mathcal{O}_{2,h} \times \dots \times \mathcal{O}_{n,h}$  to denote the joint observation of all the  $n$  agents at timestep  $h$ . We use  $\mathbb{O} = \{\mathbb{O}_h\}_{h \in [H]}$  to denote the collection of emission functions, where  $o_h \sim \mathbb{O}_h(\cdot | s_h) \in \Delta(\mathcal{O}_h)$  at any state  $s_h \in \mathcal{S}$ . We denote by  $\mathbb{O}_{i,h}(\cdot | s_h)$  the emission for agent  $i$ , which is the marginal distribution of  $o_{i,h}$  given  $\mathbb{O}_h(\cdot | s_h)$ , at any  $s_h \in \mathcal{S}$ . At each timestep  $h$ , agents will receive a common reward  $r_h = \mathcal{R}_h(s_h, a_h)$ , where  $\mathcal{R}_h : \mathcal{S} \times \mathcal{A}_h \rightarrow [0, 1]$  is a common reward function shared by the agents.

2) *Communication components*: In addition to reward-driven decision-making, agents also need to decide and learn (what) to communicate with others. At timestep  $h \in [H]$ , agents share part of their information  $z_h \in \mathcal{Z}_h$  with other agents, where  $z_h$  consists of two parts, the *baseline-sharing* part  $z_h^b \in \mathcal{Z}_h^b$  that comes from some existing sharing protocol among agents, and the *additional-sharing* part  $z_h^a \in \mathcal{Z}_h^a$  for each agent  $i$  that comes from explicit communication to be *decided/learned*, with the joint additional-sharing information  $z_h^a := \cup_{i=1}^n z_{i,h}^a$ . We keep the baseline sharing (which may be void) for generality, since a certain amount of sharing is necessary for computational and sample tractability [15], the focus of this paper. Note that  $z_h = z_h^b \cup z_h^a$  and  $z_h^b \cap z_h^a = \emptyset$ . The shared information is part of the historical observations and (both *control* and *communication*) actions.

At timestep  $h$ , the *common information* among all the agents is thus defined as:  $c_{h-} = \cup_{t=1}^{h-1} z_t \cup z_h^b$ , and  $c_{h+} = \cup_{t=1}^h z_t$ , where  $c_{h-}$  and  $c_{h+}$  denote the (accumulated) common information *before* and *after* additional sharing, respectively. The *private information* of agent  $i$  at time  $h$  *before* and *after* additional sharing are denoted by  $p_{i,h-}, p_{i,h+}$ , respectively, where  $p_{i,h-} \subseteq \{o_{i,1}, a_{i,1}, \dots, a_{i,h-1}, o_{i,h}\} \setminus c_{h-}$ ,  $p_{i,h+} \subseteq \{o_{i,1}, a_{i,1}, \dots, a_{i,h-1}, o_{i,h}\} \setminus c_{h+}$ . We denote by  $p_{h-} := (p_{1,h-}, \dots, p_{n,h-})$  and  $p_{h+} := (p_{1,h+}, \dots, p_{n,h+})$  the joint private information *before* and *after* additional sharing. We then denote by  $\tau_{i,h-} := p_{i,h-} \cup c_{h-}$ ,  $\tau_{i,h+} := p_{i,h+} \cup c_{h+}$  the *information available* to agent  $i$  at timestep  $h$ , *before* and *after* additional sharing, with  $\tau_{h-} := p_{h-} \cup c_{h-}$ ,  $\tau_{h+} := p_{h+} \cup c_{h+}$  denoting the associated *joint information*.

We use  $m_{i,h} \in \mathcal{M}_{i,h}$  to denote the *communication action* of agent  $i$  at timestep  $h$ , determining what information  $z_{i,h}^a$  she will share, through the way to be specified later. We denote by  $m_h := (m_{1,h}, \dots, m_{n,h}) \in \mathcal{M}_h := \prod_{i \in [n]} \mathcal{M}_{i,h}$  the joint communication action of all the agents. We use  $\mathcal{K}_h : \mathcal{Z}_h^a \rightarrow [0, 1]$  to denote the *communication cost* function, and  $\kappa_h = \mathcal{K}_h(z_h^a)$  to denote the communication cost at timestep  $h$ . The information flow evolves as follows, where we follow

the convention that any quantity at  $h = 0$  is empty/null.

### Assumption II.1 (Information evolution).

- (a) (Baseline sharing). For each  $h \in [H]$ ,  $z_h^b = \chi_h(p_{(h-1)+}, a_{h-1}, o_h)$  for some fixed transformation  $\chi_h$ ; for each agent  $i \in [n]$ ,  $p_{i,h-} = \xi_{i,h}(p_{i,(h-1)+}, a_{i,h-1}, o_{i,h})$  for some fixed transformation  $\xi_{i,h}$ , and the joint private information thus evolves as  $p_{h-} = \xi_h(p_{(h-1)+}, a_{h-1}, o_h)$  for some fixed transformation  $\xi_h$ ;
- (b) (Additional sharing). For each  $i \in [n], h \in [H]$ ,  $z_{i,h}^a = \phi_{i,h}(p_{i,h-}, m_{i,h})$  for some function  $\phi_{i,h}$ , given communication action  $m_{i,h}$ , and moreover,  $m_{i,h} \in z_{i,h}^a$ ; the joint additional sharing information  $z_h^a := \cup_{i \in [n]} z_{i,h}^a$  is thus generated by  $z_h^a = \phi_h(p_{h-}, m_h)$ , for some function  $\phi_h$ ; for each agent  $i \in [n]$ ,  $p_{i,h+} = p_{i,h-} \setminus z_{i,h}^a$ ;
- (c) ( $(\tau_{i,h-}, \tau_{i,h+})$ -inclusion). For each  $i \in [n], h \in [H]$ ,  $\tau_{i,h-} \subseteq \tau_{i,h+} \subseteq \tau_{i,(h+1)-}$ , and  $o_{i,h} \in \tau_{i,h-}$ .

Note that the *fixed transformations* above (e.g., the  $\chi_h$  and  $\xi_{i,h}$ ) are not affected by the *realized values* of the random variables, but dictate some *pre-defined* transformation of the input random variables. See [13], [14] and §B in [15] for common examples of baseline sharing that admit such fixed transformations, and examples in [1, §A] on how they are extended to the LTC setting. Condition (c) above assumes that the agent has full memory of the information she had in the past and at present. We emphasize that this is closely related, but different from the common notion of *perfect recall* [16], where the agent has to also recall *all her past actions*. Condition (c), in contrast, relaxes the memorization of the actions, but includes the *instantaneous observation*  $o_{i,h}$ , as a basic requirement for *closed-loop* decision-making/control. This condition is satisfied by the models and examples in [10], [13], [14], [15], and see also [1, §A] for more examples that satisfy this assumption.

Meanwhile, for both the baseline and additional sharing protocols, we follow the model in the series of studies on partial history/information sharing [13], [14], [15], [7], [8] that, if an agent shares, she will share the information with *all other* agents. We make it formal below using the verbiage with  $\sigma$ -algebra, in order to be compatible with the literature on information structures [17], [10] to be discussed later.

**Assumption II.2.**  $\forall i_1, i_2 \in [n], h_1, h_2 \in [H], i_1 \neq i_2, h_1 < h_2$ , if  $\sigma(o_{i_1, h_1}) \subseteq \sigma(\tau_{i_2, h_2-})$ , then  $\sigma(o_{i_1, h_1}) \subseteq \sigma(c_{h_2-})$ , and if  $\sigma(a_{i_1, h_1}) \subseteq \sigma(\tau_{i_2, h_2-})$ , then  $\sigma(a_{i_1, h_1}) \subseteq \sigma(c_{h_2-})$ ; if  $\sigma(o_{i_1, h_1}) \subseteq \sigma(\tau_{i_2, h_2+})$ , then  $\sigma(o_{i_1, h_1}) \subseteq \sigma(c_{h_2+})$ , and if  $\sigma(a_{i_1, h_1}) \subseteq \sigma(\tau_{i_2, h_2+})$ , then  $\sigma(a_{i_1, h_1}) \subseteq \sigma(c_{h_2+})$ .

Assumptions II.1-II.2 will be made throughout the paper.

3) *Strategies and solution concept*: At timestep  $h$ , each agent  $i$  has two strategies, a *control* strategy and a *communication* strategy. We define a control strategy as  $g_{i,h}^a : \tau_{i,h+} \rightarrow \mathcal{A}_{i,h}$  and a communication strategy as  $g_{i,h}^m : \tau_{i,h-} \rightarrow \mathcal{M}_{i,h}$ . We denote by  $g_h^a = (g_{1,h}^a, \dots, g_{n,h}^a)$  the joint control strategy and by  $g_h^m = (g_{1,h}^m, \dots, g_{n,h}^m)$  the joint communication strategy. We denote by  $\mathcal{G}_{i,h}^a, \mathcal{G}_{i,h}^m, \mathcal{G}_h^a, \mathcal{G}_h^m$  the corresponding spaces of  $g_{i,h}^a, g_{i,h}^m, g_h^a, g_h^m$ , respectively.

The objective of the agents in the LTC problem is to maximize the expected accumulated sum of the reward and the negative communication cost from timestep  $h = 1$  to  $H$ :

$$J_{\mathcal{L}}(g_{1:H}^a, g_{1:H}^m) := \mathbb{E}_{\mathcal{L}} \left[ \sum_{h=1}^H (r_h - \kappa_h) \mid g_{1:H}^a, g_{1:H}^m \right],$$

where the expectation  $\mathbb{E}_{\mathcal{L}}$  is taken over all the randomness in  $\mathcal{L}$ , given the strategies  $(g_{1:H}^a, g_{1:H}^m)$ . With this objective, for any  $\epsilon \geq 0$ , we can define the solution concept of an  $\epsilon$ -team optimum for  $\mathcal{L}$  as follows.

**Definition II.3** ( $\epsilon$ -team optimum). We call a joint strategy  $(g_{1:H}^a, g_{1:H}^m)$  an  $\epsilon$ -team optimal strategy of the LTC  $\mathcal{L}$  if

$$\max_{\tilde{g}_{1:H}^a \in \mathcal{G}_{1:H}^a, \tilde{g}_{1:H}^m \in \mathcal{G}_{1:H}^m} J_{\mathcal{L}}(\tilde{g}_{1:H}^a, \tilde{g}_{1:H}^m) - J_{\mathcal{L}}(g_{1:H}^a, g_{1:H}^m) \leq \epsilon.$$

If  $\epsilon = 0$ , we call  $(g_{1:H}^a, g_{1:H}^m)$  a team-optimal strategy of  $\mathcal{L}$ .

### B. Information Structure of LTC

In decentralized stochastic control, the notion of information structure [17], [10] captures *who knows what and when* as the system evolves. In LTC, as the additional sharing via communication will also affect the IS and is *not* determined *beforehand*, when we discuss the *IS of an LTC problem*, we will refer to that of the problem *with only baseline sharing*. In particular, an LTC  $\mathcal{L}$  *without additional sharing* is essentially a Dec-POMDP (with potential baseline information sharing), and will be referred to as the *Dec-POMDP induced by  $\mathcal{L}$*  (see a formal definition in [1, §II-B] for completeness).

In §II-A, we introduced LTC in the *state-space model*. Information structure, meanwhile, is usually more conveniently discussed within the equivalent framework of the *intrinsic model* [17]. In an intrinsic model, each agent only *acts once* throughout the system evolution, and the same agent in the state-space model at different timesteps is now treated as *different agents*. There are thus  $n \times H$  agents in total. Formally, for completeness, we extend the intrinsic-model-based reformulation of LTCs in [1, §F].

(Strictly) quasi-classical ISs are important subclasses of ISs, which have been extensively studied in stochastic control [17], [11] (see the instantiation for Dec-POMDPs in [1, §F]). We extend such a categorization to LTC problems with different ISs (of the baseline sharing) as follows.

**Definition II.4** ((Strictly) quasi-classical LTC). We call an LTC  $\mathcal{L}$  (strictly) *quasi-classical* if the *Dec-POMDP induced by  $\mathcal{L}$*  (denoted by  $\overline{\mathcal{D}}_{\mathcal{L}}$ ), i.e., LTC without additional sharing, is (strictly) *quasi-classical*. Namely, each agent in the intrinsic model of  $\overline{\mathcal{D}}_{\mathcal{L}}$  knows the information (and the actions) of the agents who influence her, either directly or indirectly.

An LTC  $\mathcal{L}$  that is not QC will thus be referred to as a *non-classical* LTC. See [1, §A] for the examples of (s)QC LTCs. Note that the categorization above is based on the ISs *before* additional sharing, an inherent property of the problem.

## III. HARDNESS AND STRUCTURAL ASSUMPTIONS

It is known that computing an (approximate) team-optimum in Dec-POMDPs, which are LTCs *without* information-sharing, is NEXP-hard [12]. The hardness cannot be fully circumvented even when agents are allowed

to share information: even if agents share all the information, the LTC problem reduces to a Partially Observable Markov Decision Process (POMDP), which is known to be PSPACE-hard [18]. Hence, additional assumptions are necessary to make LTCs computationally tractable. We introduce several such assumptions and their justifications below, whose proofs can all be found in [1, §B].

Recently, [19] showed that *observable* POMDPs [20], a class of POMDPs with relatively *informative* observations, admit *quasi-polynomial time* algorithms to solve. Such a condition and quasi-polynomial complexity result was then established for Dec-POMDPs with information sharing in [15]. As solving LTCs is at least as hard as solving the Dec-POMDPs considered in [15], we first also make such an observability assumption on the *joint* emission function as in [15], to avoid computationally intractable oracles.

**Assumption III.1** ( $\gamma$ -observability [20], [19], [15]). There exists a  $\gamma > 0$  such that  $\forall h \in [H]$ , the emission  $\mathbb{O}_h$  satisfies that  $\forall b_1, b_2 \in \Delta(\mathcal{S})$ ,  $\|\mathbb{O}_h^\top b_1 - \mathbb{O}_h^\top b_2\|_1 \geq \gamma \|b_1 - b_2\|_1$ .

However, we show next that, Assumption III.1 is not enough when it comes to LTC, if the baseline sharing IS is not favorable, and in particular, *non-classical* [10]. The hardness persists even under a few additional assumptions to be introduced later that will make LTCs more tractable.

**Lemma III.2** (Non-classical LTCs are hard). For non-classical LTCs under Assumptions III.1, III.4, III.5, and III.7, finding an  $\frac{\epsilon}{H}$ -team optimum is PSPACE-hard.

Note that the hardness comes from the intuition that, when communication costs are high, the additional sharing from LTC will be limited, preventing the upgrade of the IS from a non-classical one to a (quasi)-classical one, which is hard with only the *joint* observability of the emission (see Assumption III.1), even along with several other assumptions.

By Lemma III.2, we will hence focus on the *quasi-classical* LTCs hereafter. Indeed, QC IS is also known to be critical for efficiently solving *linear* decentralized control [21], [22]. However, quasi-classicality may not be sufficient for LTCs, since the additional sharing may *break* the QC IS, and introduce computational hardness, as argued below.

Firstly, the breaking may result from the *communication strategies*. In particular, the general communication strategy space in §II-A.3 allows the dependence on agents' *private information*, which introduces incentives for *signaling* [10] and can also cause computational hardness, as shown next.

**Lemma III.3** (QC LTCs with full-history-dependent communication strategies are hard). For QC LTCs under Assumption III.1, together with Assumptions III.5 and III.7, computing a team-optimum in the general space of  $(\mathcal{G}_{1:H}^a, \mathcal{G}_{1:H}^m)$  with  $\mathcal{G}_{i,h}^m := \{g_{i,h}^m : \mathcal{T}_{i,h} \rightarrow \mathcal{M}_{i,h}\}$  is NP-hard.

The hardness in Lemma III.3 originates from the fact that when depending on the private/local information, determining the communication action can be made as a *Team Decision problem* (TDP) [23], which is known to be hard. This will be the case even when the instantaneous observations

are relatively observable (see Assumptions III.1-III.7).

To avoid this hardness, we thus focus on communication strategies that only condition on the *common information*. Note that, this assumption does not lose the generality in the sense that the private information  $p_{i,h-}$  can still be shared. It only means that the communication action is not *determined* by  $p_{i,h-}$ , and the additional sharing is still dictated by  $z_{i,h}^a = \phi_{i,h}(p_{i,h-}, m_{i,h})$  (see Assumption II.1), depending on  $p_{i,h-}$ .

**Assumption III.4** (Common-information-based communication strategy). The communication strategies take *common information* as input, with the following form:

$$\forall i \in [n], h \in [H], \quad g_{i,h}^m : \mathcal{C}_{h-} \rightarrow \mathcal{M}_{i,h}. \quad (\text{III.1})$$

Secondly, the breaking of QC may result from the *control strategies*: if some agent did *not* influence others in the baseline sharing (and thus these other agents did *not* have to access the agent's available information, while still satisfying QC), while she starts to influence others by *sharing* her (*useless*) control actions, this will make her *control strategies* relevant. We make the following two assumptions to avoid the related pessimistic cases, each followed by a computational hardness result when the condition is missing.

**Assumption III.5** (Control-useless action is not used).  $\forall i \in [n], h \in [H]$ , if agent  $i$ 's action  $a_{i,h}$  does not influence the state  $s_{h+1}$ , namely,  $\forall s_h \in \mathcal{S}, a_h \in \mathcal{A}_h, a'_{i,h} \in \mathcal{A}_{i,h}, a'_{i,h} \neq a_{i,h}, \mathbb{T}_h(\cdot | s_h, a_h) = \mathbb{T}_h(\cdot | s_h, (a'_{i,h}, a_{-i,h}))$ . Then,  $\forall h' > h$ , the random variable  $a_{i,h} \notin \tau_{h'-}$  and  $a_{i,h} \notin \tau_{h'+}$ .

**Lemma III.6** (QC LTCs without Assumption III.5 are hard). For QC LTCs under Assumptions III.1, III.4, and III.7, finding a team-optimum is still NP-hard.

Note that Assumption III.5 was *implicitly* made in the literature [14], [15] when there exist agents who cannot *control* the transition dynamics.

**Assumption III.7** (Other agents' emissions are non-degenerate).  $\forall h \in [H], i \in [n]$ ,  $\mathbb{O}_{-i,h}$  satisfies that  $\forall b_1, b_2 \in \Delta(\mathcal{S})$  such that  $b_1 \neq b_2$ ,  $\mathbb{O}_{-i,h}^\top b_1 \neq \mathbb{O}_{-i,h}^\top b_2$ .

**Lemma III.8** (QC LTCs without Assumption III.7 are hard). For QC LTCs under Assumptions III.1, III.4, and III.5, finding an  $\epsilon/H$ -team optimum is still PSPACE-hard.

We have justified the above assumptions by showing that missing one of them may cause computational intractability in general. Hence, Assumptions III.4, III.5, and III.7 will be made throughout the paper, unless otherwise noted. More importantly, as we will show later, as another justification, LTCs under Assumptions III.4, III.5, and III.7 can indeed *preserve* the (s)QC information structures *after* additional sharing, making it possible for the overall LTC problem to be computationally tractable. Examples that satisfy these assumptions can be found in [1, §A].

#### IV. SOLVING LTC PROBLEMS PROVABLY

We now study how to solve LTC problems provably, via either *planning* (with model knowledge) or *learning* (without model knowledge). The pipeline of our solution is shown in Fig. 1, and proofs of the results can be found in [1, §C].

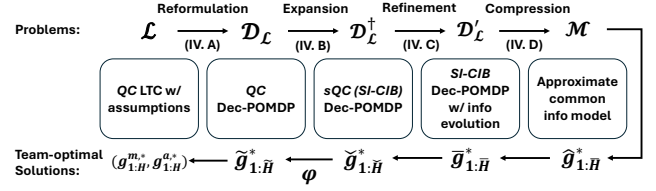


Fig. 1: The algorithmic pipeline for solving the LTC problems.

##### A. An Equivalent Dec-POMDP $\mathcal{D}_{\mathcal{L}}$

Given any LTC  $\mathcal{L}$ , we can define a Dec-POMDP  $\mathcal{D}_{\mathcal{L}}$  characterized by  $\langle \tilde{H}, \tilde{\mathcal{S}}, \{\tilde{\mathcal{A}}_{i,h}\}_{i \in [n], h \in [\tilde{H}]}, \{\tilde{\mathcal{O}}_{i,h}\}_{i \in [n], h \in [\tilde{H}]}, \{\tilde{\mathbb{T}}_h\}_{h \in [\tilde{H}]}, \{\tilde{\mathbb{O}}_h\}_{h \in [\tilde{H}]}, \tilde{\mu}_1, \{\tilde{\mathcal{R}}_h\}_{h \in [\tilde{H}]}\rangle$ , such that these two are equivalent (under the assumptions in §III):  $\forall h \in [H]$ ,

$$\begin{aligned} \tilde{H} &= 2H, \quad \tilde{\mathcal{S}} = \mathcal{S}, \quad \tilde{s}_{2h-1} = \tilde{s}_{2h} = s_h, \quad \tilde{\mathcal{A}}_{i,2h-1} = \mathcal{M}_{i,h}, \\ \tilde{\mathcal{A}}_{i,2h} &= \mathcal{A}_{i,h}, \quad \tilde{\mathcal{O}}_{i,2h-1} = \mathcal{O}_{i,h}, \quad \tilde{\mathcal{O}}_{i,2h} = \{\emptyset\}, \quad \tilde{\mu}_1 = \mu_1, \\ \tilde{\mathbb{O}}_{2h-1} &= \mathbb{O}_h, \quad \tilde{\mathbb{T}}_{2h-1}(\tilde{s}_{2h} | \tilde{s}_{2h-1}, \tilde{a}_{2h-1}) = \mathbb{1}[\tilde{s}_{2h} = \tilde{s}_{2h-1}], \\ \tilde{\mathbb{T}}_{2h}(\tilde{s}_{2h+1} | \tilde{s}_{2h}, \tilde{a}_{2h}) &= \mathbb{T}_h(\tilde{s}_{2h+1} | \tilde{s}_{2h}, \tilde{a}_{2h}), \\ \tilde{\mathcal{R}}_{2h-1} &= -\mathcal{K}_h, \quad \tilde{\mathcal{R}}_{2h} = \mathcal{R}_h, \quad \tilde{p}_{i,2h-1} = \emptyset, \quad \tilde{p}_{i,2h} = p_{i,h+}, \\ \tilde{c}_{2h-1} &= c_{h-}, \quad \tilde{c}_{2h} = c_{h+}, \quad \tilde{z}_{2h-1} = z_{h-}^b, \quad \tilde{z}_{2h} = z_h^a. \end{aligned} \quad (\text{IV.1})$$

Note that we follow the convention of  $\tilde{\tau}_{i,h} := \tilde{p}_{i,h} \cup \tilde{c}_h$  for any  $h \in [\tilde{H}]$ , and at the odd timestep  $2t-1$  for any  $t \in [H]$ , we have  $\tilde{p}_{i,2t-1} = \emptyset$  under Assumption III.4, i.e., in  $\mathcal{D}_{\mathcal{L}}$ , each agent only uses the common information so far for decision-making at timestep  $2t-1$ . Correspondingly, for any  $h \in [\tilde{H}], i \in [n]$ , we denote by  $\tilde{g}_{i,h}, \tilde{g}_h$  the strategies and by  $\tilde{\mathcal{G}}_{i,h}, \tilde{\mathcal{G}}_h$  the associated strategy spaces in  $\mathcal{D}_{\mathcal{L}}$ .

Essentially, this reformulation splits the  $H$ -step decision-making and communication procedure into a  $2H$ -step one. A similar splitting of the timesteps was also used in [7], [8]. In comparison, we consider a more general setting, where the state is not decoupled, and agents are allowed to share the observations and actions at the *previous* timesteps, due to the generality of our LTC formulation. One can verify that the  $\mathcal{L}$  and  $\mathcal{D}_{\mathcal{L}}$  are equivalent in terms of solution strategies, and  $\mathcal{D}_{\mathcal{L}}$  *preserves* the QC IS from  $\mathcal{L}$ , under Assumptions III.4, III.5, and III.7 (see [1, §IV.A] for more details).

##### B. Strict Expansion of $\mathcal{D}_{\mathcal{L}}$

However, being QC does not necessarily imply  $\mathcal{D}_{\mathcal{L}}$  can be solved *without* computationally intractable oracles. Note that this is different from the continuous-space, linear quadratic case, where QC problems can be reformulated and solved efficiently [21], [22]. With discrete spaces, the recent result [15] established a concrete *quasi-polynomial-time* complexity for planning, under the *strategy independence* assumption [14] on the common-information-based beliefs [13], [14]. This SI-CIB assumption was shown critical for *computational tractability* [15] – it eliminates the need to *enumerate* the past strategies in dynamic programming, which would otherwise be prohibitively large. Thus, we need to connect QC IS to the SI-CIB condition for computational tractability.

To this end, one can first *expand* the QC  $\mathcal{D}_{\mathcal{L}}$  to a *strictly* QC problem  $\mathcal{D}_{\mathcal{L}}^\dagger$  (notation in which will have  $\dagger$ ) by adding the



actions of the agents who influence the later agents in the intrinsic model of  $\mathcal{D}_{\mathcal{L}}$  to the shared common information. One can show that this *expansion* does not change the optimal value, and the (approximate) optimal strategy of  $\mathcal{D}_{\mathcal{L}}^{\dagger}$  can be reduced to that of  $\mathcal{D}_{\mathcal{L}}$  efficiently. Note that by definition,  $\mathcal{D}_{\mathcal{L}}^{\dagger}$  preserves the (s)QC IS of  $\mathcal{D}_{\mathcal{L}}$ . More importantly, a benefit of having a *strictly* QC  $\mathcal{D}_{\mathcal{L}}^{\dagger}$  is that, it has SI-CIBs under the assumptions in §III, making it possible to be solved without computationally intractable oracles as in [15]. See more details on this expansion in [1, §IV.B].

### C. Refinement of $\mathcal{D}_{\mathcal{L}}^{\dagger}$

Despite having SI-CIBs,  $\mathcal{D}_{\mathcal{L}}^{\dagger}$  is still not eligible for applying the results in [15]: the information evolution rules of  $\mathcal{D}_{\mathcal{L}}^{\dagger}$  break those in [14], [15]. Specifically, due to Assumption III.4, we set  $\tilde{\tau}_{i,2t-1} = \tilde{c}_{2t-1}, \tilde{p}_{i,2t-1} = \emptyset, \forall t \in [H], i \in [n]$  in  $\mathcal{D}_{\mathcal{L}}$ , which violates Assumption 1 in [14], [15]. To address this issue, we propose to further *refine*  $\mathcal{D}_{\mathcal{L}}^{\dagger}$  to obtain a Dec-POMDP  $\mathcal{D}'_{\mathcal{L}}$ , which satisfies the information evolution rules. The elements in  $\mathcal{D}'_{\mathcal{L}}$  (represented with the  $\bar{\cdot}$  notation) remain the same as those in  $\mathcal{D}_{\mathcal{L}}^{\dagger}$ , except that the private information at odd steps is now refined as  $\bar{p}_{i,2t-1} = p_{i,t-}$ , and we define  $\bar{\tau}_{i,2t-1} := \bar{p}_{i,2t-1} \cup \bar{c}_{2t-1}$  for any  $t \in [H]$ .

The new Dec-POMDP  $\mathcal{D}'_{\mathcal{L}}$  is not equivalent to  $\mathcal{D}_{\mathcal{L}}^{\dagger}$  in general, since it enlarges the strategy space at the odd timesteps. However, if we define new strategy spaces in  $\mathcal{D}'_{\mathcal{L}}$  as  $\bar{\mathcal{G}}_{i,2t-1} : \bar{\mathcal{C}}_{2t-1} \rightarrow \bar{\mathcal{A}}_{i,2t-1}, \bar{\mathcal{G}}_{i,2t} : \bar{\mathcal{T}}_{i,2t} \rightarrow \bar{\mathcal{A}}_{i,2t}$  for each  $t \in [H], i \in [n]$ , and thus define  $\bar{\mathcal{G}}_h$  to be the associated joint strategy space, then solving  $\mathcal{D}'_{\mathcal{L}}$  is equivalent to finding a *best-in-class* team-optimal strategy of  $\mathcal{D}'_{\mathcal{L}}$  within the space  $\bar{\mathcal{G}}_{1:\bar{H}}$ .

**Theorem IV.1.** Let  $\mathcal{D}_{\mathcal{L}}^{\dagger}$  be an sQC Dec-POMDP generated from  $\mathcal{L}$  after reformulation and strict expansion, and  $\mathcal{D}'_{\mathcal{L}}$  be the refinement of  $\mathcal{D}_{\mathcal{L}}^{\dagger}$  as introduced above. Then, finding the optimal strategy in  $\mathcal{D}_{\mathcal{L}}^{\dagger}$  is equivalent to finding the optimal strategy of  $\mathcal{D}'_{\mathcal{L}}$  in the space  $\bar{\mathcal{G}}_{1:\bar{H}}$ , and  $\mathcal{D}'_{\mathcal{L}}$  satisfies the following information evolution rules: for each  $h \in [\bar{H}]$ :

$$\bar{c}_h = \bar{c}_{h-1} \cup \bar{z}_h, \quad \bar{z}_h = \bar{\chi}_h(\bar{p}_{h-1}, \bar{a}_{h-1}, \bar{o}_h) \\ \text{for each } i \in [n], \quad \bar{p}_{i,h} = \bar{\xi}_{i,h}(\bar{p}_{i,h-1}, \bar{a}_{i,h-1}, \bar{o}_{i,h}),$$

with some functions  $\{\bar{\chi}_h\}_{h \in [\bar{H}]}, \{\bar{\xi}_{i,h}\}_{i \in [n], h \in [\bar{H}]}$ . Furthermore, if Assumptions III.4, III.5 and III.7 hold, then  $\mathcal{D}'_{\mathcal{L}}$  has SI-CIBs with respect to the strategy space  $\bar{\mathcal{G}}_{1:\bar{H}}$ , i.e.,  $\forall h \in [\bar{H}], \bar{s}_h \in \bar{\mathcal{S}}, \bar{p}_h \in \bar{\mathcal{P}}_h, \bar{c}_h \in \bar{\mathcal{C}}_h, \bar{g}_{1:h-1}, \bar{g}'_{1:h-1} \in \bar{\mathcal{G}}_{1:h-1}$  such that  $\bar{c}_h$  is reachable under both  $\bar{g}_{1:h-1}$  and  $\bar{g}'_{1:h-1}$ :

$$\mathbb{P}_h^{\mathcal{D}'_{\mathcal{L}}}(\bar{s}_h, \bar{p}_h | \bar{c}_h, \bar{g}_{1:h-1}) = \mathbb{P}_h^{\mathcal{D}_{\mathcal{L}}^{\dagger}}(\bar{s}_h, \bar{p}_h | \bar{c}_h, \bar{g}'_{1:h-1}). \quad (\text{IV.2})$$

### D. Planning in QC LTC with Finite-Time Complexity

Now we focus on how to solve the SI-CIB Dec-POMDP  $\mathcal{D}'_{\mathcal{L}}$  without computationally intractable oracles, which has been studied in [15]. Given a Dec-POMDP  $\mathcal{D}'_{\mathcal{L}}$  with SI-CIBs, [15] proposed to construct an *expected-approximate-common-information model*  $\mathcal{M}$  through *finite memory* (as defined in [1, §C]), when  $\mathcal{D}'_{\mathcal{L}}$  is  $\gamma$ -observable. However, the Dec-POMDP  $\mathcal{D}'_{\mathcal{L}}$  obtained from LTC has two key differences from the general ones considered in [15]. First,  $\mathcal{D}'_{\mathcal{L}}$  does

not satisfy the  $\gamma$ -observability assumption *throughout* the whole  $2H$  timesteps. Fortunately, since the emissions at odd steps are still  $\gamma$ -observable, while those at even steps are unimportant as the states remain *unchanged* from the previous step, similar results of *belief contraction* and near-optimality of finite-memory truncation as in [15] can still be obtained. Second, the rewards at the odd steps can now depend on the *private information*  $\bar{p}_h$ , instead of the state  $\bar{s}_h$ . Thanks to the definition of the approximate common-information-based beliefs  $\{\mathbb{P}_h^{\mathcal{M}}(\bar{s}_h, \bar{p}_h | \hat{c}_h)\}_{h \in [H]}$  (with  $\hat{c}_h$  denoting the approximate common information compressed from  $\bar{c}_h$ ), which is the *joint* probability of  $\bar{s}_h$  and  $\bar{p}_h$ , one can still properly evaluate the rewards at the odd steps in the algorithms of [15]. Hence, we can leverage the approaches in [15] to find an approximately optimal strategy  $\bar{g}_{1:\bar{H}}^*$  through backward induction from timesteps  $h = \bar{H}$  to 1.

Note that in each step of the backward induction, a *Team Decision problem* [23] needs to be solved for each  $\hat{c}_h$ , which is known to be NP-hard in general [23]:

$$(\bar{g}_{1,h}^*(\cdot | \hat{c}_h, \cdot), \dots, \bar{g}_{n,h}^*(\cdot | \hat{c}_h, \cdot)) \leftarrow \underset{\gamma_h}{\operatorname{argmax}} Q_h^{*,\mathcal{M}}(\hat{c}_h, \gamma_h), \quad (\text{IV.3})$$

where the  $Q$ -value function and the prescription  $\gamma_h$  are defined in [1, §C]. Hence, to obtain overall computational tractability, we make the following assumption, as in [15].

**Assumption IV.2** (One-step tractability of  $\mathcal{M}$ ). The one-step Team Decision problems induced by  $\mathcal{M}$  (i.e., Eq. (IV.3)) can be solved in polynomial time for all  $h = 2t, t \in [H]$ .

Several remarks are in order regarding Assumption IV.2. First, it can be viewed as a *minimal* assumption when it comes to computational tractability – even with  $H = 1$  and no LTC, one-step TDP requires additional structures to be solved efficiently. Second, since the Dec-POMDP here is reformulated from an LTC under Assumption III.4, it suffices to only assume one-step tractability for the *control* (i.e., even) steps. Third, even without Assumption IV.2, the SI-CIB property of  $\mathcal{D}'_{\mathcal{L}}$  and thus the derivation of dynamic programs of *fixed, tractable sizes* to solve  $\mathcal{L}$  efficiently still hold. Without such efforts, intractably many TDPs may need to be solved, leaving it less hopeful for computational tractability (even under Assumption IV.2). Finally, such an assumption is satisfied for several classes of Dec-POMDPs with information sharing, see [1, §G] for more examples. With this assumption, we can obtain a planning algorithm with quasi-polynomial time complexity (see [1, §C]).

### E. LTC with Finite-Time and Sample Complexities

Based on the planning results, we are now ready to solve the *learning* problem with both time and sample complexity guarantees. In particular, we can treat the samples from  $\mathcal{L}$  as the samples from  $\mathcal{D}'_{\mathcal{L}}$ : the *reformulation* step (§IV-A) does not change the system dynamics, but only maps the information to different random variables; the *expansion* step (§IV-B) only requires agents to share more actions with each other, without changing the input and output of the environment; the *refinement* step (§IV-C) only recovers the

private information the agents had in the original  $\mathcal{L}$ . This way, we can utilize similar algorithmic ideas in [15] to develop a learning algorithm for LTC problems. See [1, §C] for more details of the provable LTC algorithms adapted from [15]. The algorithm has the following finite-time and sample complexity guarantees.

**Theorem IV.3.** Given any QC LTC problem  $\mathcal{L}$  satisfying Assumptions III.1, III.4, III.5, and III.7, we can construct an SI-CIB Dec-POMDP problem  $\mathcal{D}'_{\mathcal{L}}$ . Moreover, there exists an LTC algorithm (see Algorithm 2 in [1, §C]) learning in  $\mathcal{D}'_{\mathcal{L}}$ , such that if the learned expected-approximate-common-information models  $\hat{\mathcal{M}}$  satisfy Assumption IV.2, then an  $\epsilon$ -team-optimal strategy for  $\mathcal{L}$  can be learned with high probability, with time and sample complexities polynomial in the parameters of  $\hat{\mathcal{M}}$ . Specifically, if  $\mathcal{L}$  has the baseline sharing protocols as in [1, §A], then such an algorithm can learn an  $\epsilon$ -team optimal strategy for  $\mathcal{L}$  with high probability, with both quasi-polynomial time and sample complexities.

## V. SOLVING GENERAL QC DEC-POMDPs

In §IV, we developed a pipeline for solving a special class of QC Dec-POMDPs generated by LTCs, without computationally intractable oracles. In fact, the pipeline can also be extended to solving general QC Dec-POMDPs, which thus advances the results in [15] that can only address SI-CIB Dec-POMDPs, a result of independent interest. Without much confusion given the context, we will adapt the notation for LTCs to studying general Dec-POMDPs: we set  $h^+ = h^- = h$  and void the additional sharing protocol. We extend the results in §IV to general QC Dec-POMDPs as follows.

**Theorem V.1.** Consider a Dec-POMDP  $\mathcal{D}$  under Assumption II.1 (c). If  $\mathcal{D}$  is sQC and satisfies Assumptions II.2, III.5, and III.7, then it has SI-CIBs. Meanwhile, if  $\mathcal{D}$  has SI-CIBs and perfect recall, then it is sQC (up to null sets).

Perfect recall [16] here means that the agents will never forget their own past information and actions (see also [1, §D]). Note that Assumption II.1 (c) is similar but different from perfect recall: it is implied by the latter with  $o_{i,h} \in \tau_{i,h}$ . Also, Assumptions II.2, III.5, and III.7 were originally made for LTCs, and here we meant to impose them for Dec-POMDPs with  $h^+ = h^- = h$ . Finally, by sQC up to null sets, we meant that if agent  $(i_1, h_1)$  influences agent  $(i_2, h_2)$  in the intrinsic model of the Dec-POMDP, then under any strategy  $\bar{g}_{1:\bar{H}}, \sigma(\bar{\tau}_{i_1, h_1}) \subseteq \sigma(\bar{\tau}_{i_2, h_2})$  except the null sets generated by  $\bar{g}_{1:\bar{H}}$ , where we add  $\bar{\cdot}$  for all the notation in the Dec-POMDP (as that of  $\mathcal{D}'_{\mathcal{L}}$  in §IV-C). Given Theorem V.1 and the results in §IV, we illustrate the relationship between LTCs and Dec-POMDPs with different assumptions and information structures in Fig. 2, which may be of independent interest.

## REFERENCES

[1] X. Liu, H. You, and K. Zhang, “Principled learning-to-communicate with quasi-classical information structures,” Tech. Rep., 2025. [Online]. Available: [https://kzhang66.github.io/assets/LTC\\_full.pdf](https://kzhang66.github.io/assets/LTC_full.pdf)

[2] J. Foerster, I. A. Assael, N. De Freitas, and S. Whiteson, “Learning to communicate with deep multi-agent reinforcement learning,” in *NeurIPS*, 2016.

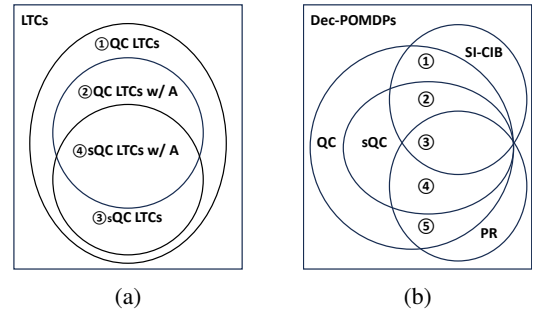


Fig. 2: (a) LTCs with different information structures. “w/ A” denotes “with Assumptions III.4, III.5, and III.7”. (b) General Dec-POMDPs with different information structures. “PR” denotes “perfect recall”. We also construct some examples for each area in [1, §H].

[3] S. Sukhbaatar, R. Fergus, et al., “Learning multiagent communication with backpropagation,” in *NeurIPS*, 2016.

[4] J. Jiang and Z. Lu, “Learning attentional communication for multi-agent cooperation,” in *NeurIPS*, 2018.

[5] S. Tatikonda and S. Mitter, “Control under communication constraints,” *IEEE Trans. Autom. Control*, vol. 49, pp. 1056–1068, 2004.

[6] S. Yüksel, “Jointly optimal LQG quantization and control policies for multi-dimensional systems,” *IEEE Trans. Autom. Control*, vol. 59, pp. 1612–1617, 2013.

[7] S. Sudhakara, D. Kartik, R. Jain, and A. Nayyar, “Optimal communication and control strategies in a multi-agent mdp problem,” *arXiv preprint arXiv:2104.10923*, 2021.

[8] D. Kartik, S. Sudhakara, R. Jain, and A. Nayyar, “Optimal communication and control strategies for a multi-agent system in the presence of an adversary,” in *IEEE Conf. on Dec. and Control*, 2022.

[9] H. S. Witsenhausen, “Separation of estimation and control for discrete time systems,” *Proceed. of the IEEE*, vol. 59, pp. 1557–1566, 1971.

[10] A. Mahajan, N. C. Martins, M. C. Rotkowitz, and S. Yüksel, “Information structures in optimal decentralized control,” in *IEEE Conf. on Dec. and Control*, 2012.

[11] S. Yüksel and T. Başar, *Stochastic Teams, Games, and Control under Information Constraints*. Springer Nature, 2023.

[12] D. S. Bernstein, R. Givan, N. Immerman, and S. Zilberstein, “The complexity of decentralized control of markov decision processes,” *Math. Oper. Res.*, vol. 27, pp. 819–840, 2002.

[13] A. Nayyar, A. Mahajan, and D. Teneketzis, “Decentralized stochastic control with partial history sharing: A common information approach,” *IEEE Trans. Autom. Control*, vol. 58, no. 7, pp. 1644–1658, 2013.

[14] A. Nayyar, A. Gupta, C. Langbort, and T. Başar, “Common information based Markov perfect equilibria for stochastic games with asymmetric information: Finite games,” *IEEE Trans. Autom. Control*, vol. 59, pp. 555–570, 2013.

[15] X. Liu and K. Zhang, “Partially observable multi-agent reinforcement learning with information sharing,” *arXiv preprint arXiv:2308.08705 (short version accepted at ICML 2023)*, 2023.

[16] H. W. Kuhn, “Extensive games and the problem of information,” in *Contrib. Theory Games, Vol. II*. Princeton Univ. Press, 1953.

[17] H. S. Witsenhausen, “The intrinsic model for discrete stochastic control: Some open problems,” in *Control Theory, Numer. Methods Comput. Syst. Model., Int. Symp., Rocquencourt*, 1975, pp. 322–335.

[18] C. H. Papadimitriou and J. N. Tsitsiklis, “The complexity of Markov decision processes,” *Math. Oper. Res.*, vol. 12, pp. 441–450, 1987.

[19] N. Golowich, A. Moitra, and D. Rohatgi, “Planning and learning in partially observable systems via filter stability,” in *Proc. 55th Annu. ACM Symp. Theory Comput.*, 2023.

[20] E. Even-Dar, S. M. Kakade, and Y. Mansour, “The value of observation for monitoring dynamic systems,” in *IJCAI*, 2007.

[21] Y.-C. Ho et al., “Team decision theory and information structures in optimal control problems – part i,” *IEEE Trans. Autom. Control*, vol. 17, pp. 15–22, 1972.

[22] A. Lamperski and L. Lessard, “Optimal decentralized state-feedback control with sparsity and delays,” *Automatica*, pp. 143–151, 2015.

[23] J. Tsitsiklis and M. Athans, “On the complexity of decentralized decision making and detection problems,” *IEEE Trans. Autom. Control*, vol. 30, pp. 440–446, 1985.