# Sequential Information Design: Learning to Persuade in the Dark

Anonymous Author(s) Affiliation Address email

### Abstract

1 We study a repeated *information design* problem faced by an informed *sender* who 2 tries to influence the behavior of a self-interested *receiver*. We consider settings where the receiver faces a sequential decision making (SDM) problem. At each 3 round, the sender observes the realizations of random events in the SDM problem. 4 This begets the challenge of how to incrementally disclose such information to 5 the receiver to *persuade* them to follow (desirable) action recommendations. We 6 study the case in which the sender does *not* know random events probabilities, and, 7 thus, they have to gradually learn them while persuading the receiver. We start by 8 providing a non-trivial polytopal approximation of the set of sender's persuasive 9 information structures. This is crucial to design efficient learning algorithms. Next, 10 we prove a negative result: no learning algorithm can be persuasive. Thus, we 11 12 relax persuasiveness requirements by focusing on algorithms that guarantee that the receiver's regret in following recommendations grows sub-linearly. In the 13 *full-feedback* setting—where the sender observes *all* random events realizations—, 14 we provide an algorithm with  $O(\sqrt{T})$  regret for both the sender and the receiver. 15 16 Instead, in the bandit-feedback setting-where the sender only observes the realizations of random events actually occurring in the SDM problem—, we design an al-17 gorithm that, given an  $\alpha \in [1/2, 1]$  as input, ensures  $\tilde{O}(T^{\alpha})$  and  $\tilde{O}(T^{\max\{\alpha, 1-\frac{\alpha}{2}\}})$ 18 19 regrets, for the sender and the receiver respectively. This result is complemented by a lower bound showing that such a regrets trade-off is essentially tight. 20

### 21 **1 Introduction**

Bayesian persuasion [19] (a.k.a. *information design*) is the problem faced by an informed *sender* who wants to influence the behavior of a self-interested *receiver* via the provision of payoff-relevant information. This captures the problem of "who gets to know what", which is fundamental in all economic interactions. Thus, Bayesian persuasion is ubiquitous in real-world problems, such as, *e.g.*, online advertising [5], voting [1, 7, 8], traffic routing [4, 10], security [23, 26], and marketing [3, 6].

We study Bayesian persuasion in settings where the receiver plays in a sequential decision making 27 (SDM) problem. An SDM problem is characterized by a tree structure made by: decision nodes, 28 where the receiver takes actions, and *chance* nodes, in which *partially observable* random events occur. 29 30 The sender perfectly observes the realizations of random events, and their goal is to incrementally disclose the acquired information to induce the receiver towards desirable outcomes. In order to 31 do so, the sender commits to a *signaling scheme* specifying a probability distribution over action 32 recommendations for the receiver at each decision node. Specifically, the sender commits to a 33 *persuasive* signaling scheme, meaning that the receiver is incentivized to follow recommendations. 34 We consider the case of a *farsighted* receiver, meaning that they take into account all the possible 35 future events when deciding whether to deviate or not from recommendations at each decision node. 36

Submitted to 36th Conference on Neural Information Processing Systems (NeurIPS 2022). Do not distribute.

With some notable exceptions (see, *e.g.*, [27]), Bayesian persuasion models in the literature make the stringent assumption that both the sender and the receiver know the *prior*, which, in our setting, is defined by the probabilities of random events in the SDM problem. We relax such an assumption by considering an online learning framework in which the sender, without any knowledge of the prior, repeatedly interacts with the receiver to gradually learn the prior while still being persuasive.

**Original contributions.** Our goal is to design online learning algorithms that are no-regret for 42 the sender, while being persuasive for the receiver. We start by providing a non-trivial polytopal 43 approximation of the set of sender's persuasive signaling schemes. This will be crucial in designing 44 efficient (*i.e.*, polynomial-time) learning algorithms, and it also shows how a sender-optimal signaling 45 scheme can be found in polynomial time in the offline version of our problem, which may be of 46 independent interest. Next, we prove a negative result: without knowing the prior, no algorithm can 47 be persuasive at each round with high probability. Thus, we relax persuasiveness requirements by 48 focusing on learning algorithms that guarantee that the receiver's regret in following recommendations 49 grows sub-linearly, while guaranteeing the same for sender's regret. First, we study the *full-feedback* 50 case, where the sender observes the realizations of *all* the random events that may potentially happen 51 in the SDM problem. In such a setting, we provide an algorithm with  $\tilde{O}(\sqrt{T})$  regret for both the 52 sender and the receiver. Then, we focus on the *bandit-feedback* setting, where the sender only 53 observes the realizations of random events on the path in the tree traversed during the SDM problem. 54 In this case, we design an algorithm that achieves  $\tilde{O}(T^{\alpha})$  sender's regret and  $\tilde{O}(T^{\max\{\alpha,1-\frac{\alpha}{2}\}})$ 55 receiver's regret, for any  $\alpha \in [1/2, 1]$  given as input. The crucial component of the algorithm is a 56 non-trivial exploration phase that uniformly explores the tree defining the SDM problem to build 57 58 suitable estimators of the prior. This is needed since, with bandit feedback, playing a signaling scheme may provide insufficient information about its persuasiveness. Finally, we provide a lower 59 bound showing that the regrets trade off achieved by our algorithm is tight for  $\alpha \in [1/2, 2/3]$ . 60

61 **Related works.** Some works addressed Bayesian persuasion in *Markov decision processes* (MDPs). Gan et al. [17] and Wu et al. [25] show how to efficiently find a sender-optimal policy when the 62 63 receiver is *myopic* (*i.e.*, it only optimizes one-step rewards) in MDPs with infinite and finite horizon, respectively. Moreover, the former assume that the environment is known, while the latter do not. 64 These works considerably differ from ours, since we assume a farsighted receiver and also model 65 partial observability of random events.<sup>1</sup> Another work close to ours is [27], which studies a (non-66 sequential) persuasion problem in which the sender and the receiver do not know the prior and interact 67 68 online. Zu et al. [27] provide a persuasive learning algorithm, while, in our model, we show that the ignorance of the prior precludes the possibility of committing to persuasive signaling schemes, 69 and, thus, we need to resort to new techniques to circumvent the issue. Finally, Celli et al. [12] study 70 Bayesian persuasion with multiple receivers interacting in an imperfect-information sequential game. 71 Differently from ours, their model adopts a different notion of persuasiveness, known as ex ante 72 persuasiveness, and it assumes that the prior is known. Other works study learning problems in which 73 74 the sender does *not* know the receivers' payoffs (but knows the prior); see, *e.g.*, [9, 11].

## 75 2 Preliminaries

### 76 2.1 Sequential decision making problems

An instance of an SDM problem is defined by a tree structure, utilities, and random events probabilities. 77 The tree structure has a set of nodes  $\mathcal{H} := \mathcal{Z} \cup \mathcal{H}_d \cup \mathcal{H}_c$ , where:  $\mathcal{Z}$  contains all the *terminal nodes* in 78 which the problem ends (corresponding to the leaves of the tree),  $\mathcal{H}_d$  is the set of *decision nodes* in 79 which the agent acts, while  $\mathcal{H}_c$  is the set of *chance nodes* where random events occur. Given any 80 non-terminal node  $h \in \mathcal{H} \setminus \mathcal{Z}$ , we let A(h) be the set of arcs outgoing from h. If  $h \in \mathcal{H}_d$ , then 81 A(h) is the set of receiver's actions available at h, while, if  $h \in \mathcal{H}_c$ , then A(h) encodes the possible 82 outcomes of the random event occurring at h. Furthermore, the utility function  $u: \mathbb{Z} \to [0, 1]$  defines 83 the agent's payoff u(z) when the problem ends in terminal node  $z \in \mathcal{Z}$ . Finally, each chance node 84  $h \in \mathcal{H}_c$  is characterized by a probability distribution  $\mu_h \in \Delta_{A(h)}$  over the possible outcomes of the 85 corresponding random event, with  $\mu_h(a)$  denoting the probability of action  $a \in A(h)$ .<sup>2</sup> 86

<sup>&</sup>lt;sup>1</sup>Gan et al. [17] also study a model with farsighted receiver, where they show that the problem of finding a sender-optimal policy is NP-hard. Thus, they do *not* provide any algorithmic result for such a model. <sup>2</sup>For a finite set X we denote with  $\Delta_X$  the set of probability distributions over X.

In an SDM problem, the agent has *imperfect information*, since they do *not* perfectly observe the 87 outcomes of random events. Thus, the set of decision nodes  $\mathcal{H}_d$  is partitioned into *information sets* 88 (infosets for short), where an infoset  $I \subseteq \mathcal{H}_d$  is a subset of decision nodes that are indistinguishable 89 for the agent. We denote the set of infosets as  $\mathcal{I}$ . For every infoset  $I \in \mathcal{I}$  and pair of nodes  $h, h' \in I$ , 90 it must be the case that A(h) = A(h') =: A(I), otherwise the agent could distinguish between 91 the two nodes. We assume that the agent has *perfect recall*, which means that they never forget 92 information once acquired. Formally, this is equivalent to assume that, for every infoset  $I \in \mathcal{I}$ , all the 93 paths from the root of the tree to a node  $h \in I$  identify the same ordered sequence of agent's actions. 94

### 95 2.2 Bayesian persuasion in sequential decision making problems

We study Bayesian persuasion in SDM (BPSDM) problems. These extend the classical Bayesian 96 persuasion framework [19] to SDM problems by introducing an exogenous agent that acts as a sender 97 by issuing signals to the decision-making agent (the *receiver*).<sup>3</sup> By following the Bayesian persuasion 98 terminology, the probability distributions  $\mu_h$  for each chance node h are collectively referred to as 99 the *prior*. Thus, the sender observes the realizations of random events occurring in the SDM problem 100 and can partially disclose information to influence the receiver's behavior. Moreover, the sender has 101 their own utility function defined over terminal nodes, denoted as  $f: \mathcal{Z} \to [0, 1]$ , and their goal is to 102 commit to a publicly known signaling scheme that maximizes their utility in expectation with respect 103 to the prior, the selected signaling scheme, and the receiver's strategy. 104

Formally, a signaling scheme for the sender defines a probability distribution  $\phi_h \in \Delta_{S(h)}$  at each decision node  $h \in \mathcal{H}_d$ , where S(h) is a finite set of signals available at h. During the SDM problem, when the receiver reaches a node  $h \in \mathcal{H}_d$  belonging to an infoset  $I \in \mathcal{I}$ , the sender draws a signal  $s \sim \phi_h$  and communicates it to the receiver. Then, based on the history of signals observed from the beginning of the SDM problem (s included), the receiver computes a *posterior* belief over the nodes belonging to the infoset I and plays so as to maximize their expected utility in the SDM sub-problem that starts from I, taking into account the just acquired information.

As customary in these settings, a simple revelation-principle-style argument allows us to focus on signaling schemes that are *direct* and *persuasive* [2, 19]. In particular, a signaling scheme is direct if signals correspond to action recommendations, namely S(h) = A(h) for all  $h \in \mathcal{H}_d$ . A direct signaling scheme is persuasive if the receiver is incentivized to follow action recommendations issued by the sender. Moreover, we assume that, if the receiver does *not* follow action recommendations at some decision node, then the sender stops issuing recommendations at nodes later reached during the SDM problem. This is without loss of generality.<sup>4</sup>

#### 119 2.3 The sequence-form representation

The *sequence form* is a commonly-used, compact way of representing (*mixed*) *strategies* in SDM problems [20]. In this work, the sequence-form representation will be employed for receiver's strategies, and to encode the signaling schemes and priors, as we describe in the following.

**Receiver's strategies.** Given any  $h \in \mathcal{H}$ , we let  $\sigma_r(h)$  be the ordered *sequence* of receiver's actions 123 on the path from the root of the tree to node h. By the perfect recall assumption, given any infoset 124  $I \in \mathcal{I}$ , it holds that  $\sigma_r(h) = \sigma_r(h') =: \sigma_r(I)$  for every pair of nodes  $h, h' \in I$ . Thus, we can identify 125 sequences with infoset-action pairs, with  $\sigma = (I, a)$  encoding the sequence of actions obtained by 126 appending action  $a \in A(I)$  at the end of  $\sigma_r(I)$ , for any infoset  $I \in \mathcal{I}$ . Moreover,  $\emptyset$  denotes the 127 *empty sequence*. Hence, the receiver's sequences are  $\Sigma_r := \{(I, a) \mid I \in \mathcal{I}, a \in A(I)\} \cup \{\emptyset\}$ . In the 128 sequence-form representation, mixed strategies are defined by specifying the probability of playing 129 each sequence of actions. Thus, a receiver's strategy is represented by a vector  $x \in [0, 1]^{|\Sigma_r|}$ , where 130  $x[\sigma]$  encodes the realization probability of sequence  $\sigma \in \Sigma_r$ . Furthermore, a sequence-form strategy 131 is well-defined if and only if it satisfies the following linear constraints: 132

$$\boldsymbol{x}[\varnothing] = 1$$
 and  $\boldsymbol{x}[\sigma_r(I)] = \sum_{a \in A(I)} \boldsymbol{x}[\sigma_r(I)a] \quad \forall I \in \mathcal{I}.$ 

We denote by  $\mathcal{X}_r$  the polytope of all receiver's sequence-form strategies. We will also need to work with the sets of receiver's strategies in the SDM sub-problem that starts from an infoset  $I \in \mathcal{I}$ , formally defined as  $\mathcal{X}_{r,I} \coloneqq \{ x \in \mathcal{X}_r \mid x[\sigma_r(I)] = 1 \}$ .

<sup>&</sup>lt;sup>3</sup>Appendix A shows that BPSDM reduces to classical Bayesian persuasion when there is no sequentiality.

<sup>&</sup>lt;sup>4</sup> For a discussion on a similar problem in the field of correlation in sequential games, we refer to [22, 24].

Signaling schemes. We represent signaling schemes in sequence form by leveraging the fact that the sender can be thought of as a perfect-information agent who plays at the decision nodes of the SDM

problem, since their actions correspond to recommendations for the receiver. Thus, since sender's

infosets correspond to decision nodes, their sequences  $\Sigma_s := \{(h, a) \mid h \in \mathcal{H}_d, a \in A(h)\} \cup \{\emptyset\}.$ 

Then, we denote the polytope of (sequence-form) signaling schemes as  $\Phi \subseteq [0, 1]^{|\Sigma_s|}$ , where each

141 signaling scheme is represented as a vector  $\boldsymbol{\phi} \in [0,1]^{|\Sigma_s|}$  satisfying:

 $\phi[\varnothing] = 1$  and  $\phi[\sigma_s(h)] = \sum_{a \in A(h)} \phi[\sigma_s(h)a] \quad \forall h \in \mathcal{H}_d,$ 

where, similarly to  $\sigma_r(h)$  for the receiver,  $\sigma_s(h)$  denotes the sender's sequence identified by  $h \in \mathcal{H}$ . We also define  $\Pi := \Phi \cap \{0, 1\}^{|\Sigma_s|}$  as the set of *deterministic* signaling schemes, which are those that recommend a single action with probability one at each decision node.

Priors. We also encode prior probability distributions  $\mu_h$  by means of the sequence form. Indeed, these can be though of as elements of a fixed strategy played by a (fictitious) perfect-information agent that acts at chance nodes. Thus, for such a chance agent, we define  $\Sigma_c$ ,  $\mathcal{X}_c$ , and  $\sigma_c(h)$  as their counterparts previously introduced for the receiver. Moreover, in the following, we denote by  $\mu^* \in \mathcal{X}_c$  the (sequence-form) prior, recursively defined as follows:

 $\boldsymbol{\mu}^{\star}[\varnothing] \coloneqq 1 \quad \text{and} \quad \boldsymbol{\mu}^{\star}[\sigma_c(h)a] \coloneqq \boldsymbol{\mu}^{\star}[\sigma_c(h)] \, \boldsymbol{\mu}_h(a) \quad \forall h \in \mathcal{H}_c, \forall a \in A(h).$ 

**Ordering of sequences.** For the sake of presentation, we introduce a partial ordering relation among sequences. Given two sequences  $\sigma = (I, a) \in \Sigma_r$  and  $\sigma' = (J, b) \in \Sigma_r$ , we write  $\sigma \preceq \sigma'$ (read as  $\sigma$  precedes  $\sigma'$ ), whenever there exists a path in the tree connecting a node in I to a node in J, and such a path includes action a. We adopt analogous definitions for sequences in  $\Sigma_s$  and  $\Sigma_c$ .<sup>5</sup>

### 154 **3** Learning to persuade

In this work, we relax the strong assumption that both the sender and the receiver know the prior  $\mu^{\star}$ 155 by casting the BPSDM problem into an *online learning framework* in which the sender repeatedly 156 interacts with the receiver over a time horizon of length T. At each round  $t \in [T]$ , the interaction 157 goes as follows:<sup>6</sup> (i) the sender commits to a signaling scheme  $\phi_t \in \Phi$ ; (ii) a vector  $y_t \in \{0,1\}^{|\Sigma_c|}$ 158 encoding realizations of random events is drawn according to  $\mu^*$ ; (iii) the sender and the receiver play 159 an instance of the (one-shot) BPSDM problem (detailed in Section 2.2), in which the sender commits 160 to  $\phi_t$ , random events at chance nodes are realized as defined by  $y_t$ , and the receiver sticks to the 161 recommendations issued by the sender; and (iv) the sender observes a feedback on the realization of 162 163 random events at chance nodes, which can be of two types: *full feedback* when the sender observes  $y_i$ , which specifies the realizations of *all* the random events at chance nodes that are possibly reachable 164 during the SDM problem; *bandit feedback* when the sender observes the terminal node  $z_t \in \mathcal{Z}$ 165 reached at the end of the SDM problem. The latter is equivalent to observing the realizations of 166 random events at the chance nodes that are actually reached during the SDM problem, namely  $\sigma_c(z_t)$ . 167

By letting  $\Phi^{\diamond}(\mu^{\star})$  be the set of persuasive signaling schemes, *i.e.*, such that the receiver is incentivized to following recommendations (a formal definition is provided in Definition 2), the goal of the sender is to select a sequence of signaling schemes, namely  $\phi_1, \ldots, \phi_T$ , which maximizes their expected utility, while guaranteeing that each signaling scheme  $\phi_t$  is persuasive, namely  $\phi_t \in \Phi^{\diamond}(\mu^{\star})$ .

We measure the performance of a sequence  $\phi_1, \ldots, \phi_T$  of signaling schemes by comparing it with an optimal (fixed) persuasive signaling scheme. Formally, given a signaling scheme  $\phi \in \Phi$ , we first define  $U(\phi, \mu^*)$ , respectively  $F(\phi, \mu^*)$ , as the expected utility achieved by the receiver, respectively the sender, whenever the former follows action recommendations. These can be expressed as linear functions of  $\phi$ , which, for any  $\mu \in \mathcal{X}_c$ , are defined as follows:

$$U(\boldsymbol{\phi}, \boldsymbol{\mu}) \coloneqq \sum_{z \in \mathcal{Z}} \boldsymbol{\mu}[\sigma_c(z)] \boldsymbol{\phi}[\sigma_s(z)] u(z), \quad F(\boldsymbol{\phi}, \boldsymbol{\mu}) \coloneqq \sum_{z \in \mathcal{Z}} \boldsymbol{\mu}[\sigma_c(z)] \boldsymbol{\phi}[\sigma_s(z)] f(z).$$

Finally, by letting  $\phi^* \in \operatorname{argmax}_{\phi \in \Phi^{\diamond}(\mu^*)} F(\phi, \mu^*)$  be an optimal (fixed) persuasive signaling scheme, the sender' performance over T rounds is measured by the (*cumulative*) sender's regret:

$$R_T := \sum_{t \in [T]} \left( F(\boldsymbol{\phi}^{\star}, \boldsymbol{\mu}^{\star}) - F(\boldsymbol{\phi}_t, \boldsymbol{\mu}^{\star}) \right)$$

<sup>&</sup>lt;sup>5</sup>We refer the reader to Appendix B for an example of SDM problem and its sets of sequences.

<sup>&</sup>lt;sup>6</sup>Throughout this work, for  $n \in \mathbb{N}$ , we denote with [n] the set  $\{1, \ldots, n\}$ .

The goal is to design learning algorithms (for the sender) which select sequences of persuasive signaling schemes such that  $R_T$  grows asymptotically sub-linearly in T, namely  $R_T = o(T)$ .

### <sup>181</sup> 4 On the characterization of persuasive signaling schemes

#### 182 4.1 A local decomposition of persuasiveness

In this section, we formally introduce the set of persuasive signaling schemes  $\Phi^{\diamond}(\mu^{\star})$  as the set of signaling schemes for which the receiver's expected utility by following recommendations is greater than the one provided by an optimal *deviation policy* (DP).<sup>7</sup> In addition, we show how to decompose any DP into components defined locally at each infoset, which will be crucial in the following Section 4.2. Intuitively, a DP for the receiver is specified by two elements: (i) a set of *deviation points* in which the DP prescribes to stop following action recommendations; and (ii) the *continuation strategies* to be adopted after deviating from recommendations.

We represent deviation points by vectors  $\boldsymbol{\omega} \in \{0, 1\}^{|\Sigma_r|}$ , which are defined so that  $\boldsymbol{\omega}[\sigma] = 1$  if and only if the DP prescribes to deviate upon observing the sequence of action recommendations  $\sigma \in \Sigma_r$ . Moreover, by leveraging the w.l.o.g. assumption that the sender stops issuing recommendations after the receiver deviated from them, we focus on DPs such that each path from the root of the tree to a terminal node involves only one deviation point. As a result, the set of all valid vectors  $\boldsymbol{\omega} \in \{0,1\}^{|\Sigma_r|}$ is formally defined as  $\Omega \coloneqq \left\{ \boldsymbol{\omega} \in \{0,1\}^{|\Sigma_r|} \mid \sum_{\sigma \in \Sigma_r: \sigma \preceq \sigma_r(z)} \boldsymbol{\omega}[\sigma] \le 1 \quad \forall z \in \mathcal{Z} \right\}$ .

We represent the continuation strategies of DPs by introducing the set of *continuation strategy profiles*, denoted as  $\mathcal{P} \coloneqq \bigotimes_{\sigma = (I,a) \in \Sigma_r} \mathcal{X}_{r,I}$ . A continuation strategy profile  $\rho \in \mathcal{P}$ , with  $\rho = (\rho_{\sigma})_{\sigma \in \Sigma_r}$ , defines a strategy  $\rho_{\sigma} \in \mathcal{X}_{r,I}$  for every receiver's sequence  $\sigma = (I, a) \in \Sigma_r$ . Intuitively,  $\rho_{\sigma}$  is the strategy for the SDM sub-problem starting from infoset *I* that is used by the receiver after deviating upon observing sequence  $\sigma \in \Sigma_r$ . As a result, any pair  $(\omega, \rho) \in \Omega \times \mathcal{P}$  specifies a valid DP; formally:

**Definition 1** (Deviation policy). *Given a vector*  $\boldsymbol{\omega} \in \Omega$  *and a profile*  $\boldsymbol{\rho} \in \mathcal{P}$ *, the*  $(\boldsymbol{\omega}, \boldsymbol{\rho})$ -DP *prescribes* to follow sender's recommendations until action *a* is recommended at infoset *I* for some sequence  $\sigma = (I, a)$  such that  $\boldsymbol{\omega}[\sigma] = 1$ ; from that point on, it prescribes to play according to strategy  $\boldsymbol{\rho}_{\sigma}$ .

We denote by  $U^{\omega \to \rho}(\phi, \mu^*)$  the receiver's expected utility obtained with a  $(\omega, \rho)$ -DP, so that we can state the following formal definition of persuasive signaling schemes.

**Definition 2** (Persuasiveness). A signaling scheme  $\phi \in \Phi$  is  $\epsilon$ -persuasive, namely  $\phi \in \Phi_{\epsilon}^{\circ}(\mu^{\star})$ , if

$$\max_{(\boldsymbol{\omega},\boldsymbol{\rho})\in\Omega\times\mathcal{P}} U^{\boldsymbol{\omega}\to\boldsymbol{\rho}}(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) \leq \epsilon.$$
(1)

208 Moreover, a signaling scheme  $\phi \in \Phi$  is persuasive, namely  $\phi \in \Phi^{\diamond}(\mu^{\star})$ , if it is 0-persuasive.

Intuitively, the above definition states that a signaling scheme is  $\epsilon$ -persuasive if the receiver's expected utility by following recommendations is at most  $\epsilon$  less than the one obtained by an optimal DP, which is a DP maximizing receiver's expected utility.

Our local decomposition of DPs is based on suitably-defined, simple deviation policies, which we call *single-point DPs* (SPDPs). These are a special case of DPs that stop following sender's action recommendations only when a specific single infoset is reached and a particular action is recommended therein. SPDPs are formally defined as follows:

**Definition 3** (Single-point deviation strategy). Given a receiver's sequence  $\sigma = (I, a) \in \Sigma_r$  and a receiver's strategy  $\rho_{\sigma} \in \mathcal{X}_{r,I}$  for the SDM sub-problem starting from infoset I, the  $(\sigma, \rho_{\sigma})$ -SPDP prescribes to follow sender's recommendations until action a is recommended at infoset I; from that point on, the strategy prescribes to play according to  $\rho_{\sigma}$ .

We denote by  $U_{\sigma \to \rho_{\sigma}}(\phi, \mu^{\star})$  the receiver's expected utility obtained by following an  $(\sigma, \rho_{\sigma})$ -SPDP.

<sup>221</sup> The following theorem provides the key result underlying our decomposition.<sup>8</sup> It shows that the dif-

ference between the utility achieved by a ( $\omega, \rho$ )-DP and that obtained by following recommendations

<sup>&</sup>lt;sup>7</sup>For ease of exposition, all the definitions and results in this section are provided for the prior  $\mu^*$ . It is straightforward to generalize them to the case of a generic  $\mu \in \mathcal{X}_c$ .

<sup>&</sup>lt;sup>8</sup>All the proofs are provided in the Appendices D, E, F, and G.

- can be decomposed into the sum over all the sequences  $\sigma \in \Sigma_r$  of analogous differences defined for the  $(\sigma, \rho_{\sigma})$ -SPDPs, where each difference is weighted by  $\omega[\sigma]$ .
- **Theorem 1.** Given a signaling scheme  $\phi \in \Phi$  and a  $(\omega, \rho)$ -DP, it holds:

$$U^{\boldsymbol{\omega} \to \boldsymbol{\rho}}(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) = \sum_{\sigma \in \Sigma_r} \boldsymbol{\omega}[\sigma] \Big( U_{\sigma \to \boldsymbol{\rho}_{\sigma}}(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) \Big).$$

### **4.2** A polytopal approximation of the set of persuasive signaling schemes

In the following, we show how to exploit Theorem 1 to provide an approximate characterization of the set  $\Phi_{\epsilon}^{\diamond}(\mu^{\star})$  using a polynomially-sized polytope. First, we state a corollary of Theorem 1 showing that persuasiveness can be bounded by suitably defined SPDPs. Formally:<sup>9</sup>

**Corollary 1.** Given a signaling scheme  $\phi \in \Phi$ , the following holds:

$$\max_{(\boldsymbol{\omega},\boldsymbol{\rho})\in\Omega\times\mathcal{P}} U^{\boldsymbol{\omega}\to\boldsymbol{\rho}}(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) \leq \sum_{\sigma=(I,a)\in\Sigma_r} \left| \max_{\boldsymbol{\rho}_{\sigma}\in\mathcal{X}_{r,I}} U_{\sigma\to\boldsymbol{\rho}_{\sigma}}(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) \right|^{\top}.$$

¬ +

By exploiting Corollary 1, we introduce the following definition of  $\epsilon$ -persuasive polytope (Lemma 1

justifies the term polytope), as the set of signaling schemes for which there is no  $(\sigma, \rho_{\sigma})$ -SPDP that

achieves a receiver's utility that exceeds by more than  $\epsilon/|\Sigma_r|$  that of following recommendations.

**Definition 4** (Persuasive polytope). *The*  $\epsilon$ -persuasive polytope *is defined as:* 

$$\Lambda_{\epsilon}(\boldsymbol{\mu}^{\star}) \coloneqq \Big\{ \boldsymbol{\phi} \in \Phi \ \Big| \ \max_{\boldsymbol{\rho}_{\sigma} \in \mathcal{X}_{r,I}} U_{\sigma \to \boldsymbol{\rho}_{\sigma}}(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) \leq \epsilon / |\Sigma_{r}| \quad \forall \sigma \in \Sigma_{r} \Big\}.$$

235 Moreover, we denote by  $\Lambda(\mu^*)$  the 0-persuasive polytope.

- As we show in the following lemma,  $\Lambda_{\epsilon}(\mu^{\star})$  is an efficiently-representable polytope.
- **Lemma 1.** The set  $\Lambda_{\epsilon}(\mu^{\star})$  can be described by means of a polynomial number of linear constraints.

The following lemma shows that the  $\epsilon$ -persuasive polytope is contained in  $\Phi_{\epsilon}^{\diamond}(\mu^{\star})$ , and that the set of persuasive signaling schemes is contained in the former. Formally:

**Lemma 2.** It is always the case that  $\Phi^{\diamond}(\mu^{\star}) \equiv \Lambda(\mu^{\star}) \subseteq \Lambda_{\epsilon}(\mu^{\star}) \subseteq \Phi_{\epsilon}^{\diamond}(\mu^{\star})$ .

Lemma 2 also implies that the polytope  $\Lambda(\mu^*)$  exactly characterizes the set of persuasive signaling schemes  $\Phi^{\diamond}(\mu^*)$ . Thus, by adding the maximization of the sender's expected utility  $F(\phi, \mu^*)$  on top of the linear constraints describing  $\Lambda(\mu^*)$ , we obtain a polynomially-sized linear program for finding an optimal sender's signaling scheme in any instance of the BPSDM problem in which  $\mu^*$  is known.

**Theorem 2.** The BPSDM problem can be solved in polynomial time when the prior  $\mu^*$  is known.

### <sup>246</sup> 5 Always being persuasive is impossible: a relaxation is needed

In this section, we prove that it is impossible to design an algorithm that returns a sequence of persuasive signaling schemes for a generic BPSDM problem. Motivated by this result, we introduce a new online learning problem that relaxes persuasiveness requirements.

<sup>250</sup> First, we provide the following impossibility result:

**Theorem 3** (Impossibility of persuasiveness). There exists a constant  $\gamma \in (0,1)$  such that no algorithm can guarantee to output a sequence  $\phi_1, \ldots, \phi_T$  of signaling schemes such that, with probability al least  $\gamma$ , all the signaling schemes  $\phi_t$  are persuasive.

Notice that this result is in contrast with what happens in the basic case of non-sequential Bayesian persuasion (see the work by Zu et al. [27]), where it is possible to design no-regret algorithms that output sequences of signaling schemes that are guaranteed to be persuasive with high probability.

ourban sedannes of vignaning senemes and me gamminee to se become the seneme by

Theorem 3 motivates the introduction of a less restrictive requirement on the signaling schemes output by a learning algorithm. In particular, we look for algorithms that output sequences  $\phi_1, \ldots, \phi_T$  of

<sup>9</sup>Given any  $x \in \mathbb{R}$ , we let  $[x]^+ := \max(x, 0)$ .

signaling schemes such that the expected utility loss incurred by the receiver by following sender's

recommendations rather than playing an optimal DP is small. To capture such a requirement, we

introduce the following definition of (*cumulative*) receiver's regret:

$$V_T \coloneqq \sum_{t \in [T]} \max_{(\boldsymbol{\omega}, \boldsymbol{\rho}) \in \Omega \times \mathcal{P}} U^{\boldsymbol{\omega} \to \boldsymbol{\rho}}(\boldsymbol{\phi}_t, \boldsymbol{\mu}^{\star}) - \sum_{t \in [T]} U(\boldsymbol{\phi}_t, \boldsymbol{\mu}^{\star})$$

Therefore our goal becomes that of designing algorithms guaranteeing that the cumulative receiver's regret grows sub-linearly in T, namely  $V_T = o(T)$ , while continuing to ensure that  $R_T = o(T)$ .

In Sections 6 and 7, we design algorithms achieving sub-linear  $V_T$  and  $R_T$  for the learning problem described in Section 3. The algorithms implement two functions: (i) SELECTSTRATEGY(), which, at each  $t \in [T]$ , draws a signaling scheme  $\phi_t \in \Phi$  on the basis of the internal state of the algorithm; and (ii) UPDATE $(o_t)$ , which modifies the internal state on the basis of the observation  $o_t$  received as feedback. Each algorithm alternates these two functions as the interaction between the sender and the receiver unfolds as described in Section 3. Specifically, under full feedback the sender observes  $y_t$ and calls UPDATE $(y_t)$ , while in the bandit feedback it observes  $z_t$  and calls UPDATE $(z_t)$ .

### 271 6 Learning with full feedback

We start by providing a learning algorithm (Algorithm 1) working with full feedback, *i.e.*, when the sender observes the realizations of *all* the possible random events.

The main idea of the algorithm is to choose signaling 275 schemes  $\phi_t$  that belong to suitable sets  $\Lambda_{\beta_t}(\widehat{\mu}_t)$  which are 276 designed to be "close" to the set  $\Phi^{\diamond}(\mu^{\star})$  of persuasive sig-277 naling schemes. At each round  $t \in [T]$ , Algorithm 1 defines 278 the desired set as follows. First, it maintains an estimate 279  $\hat{\mu}_t$  of  $\mu^*$ ; formally, it defines a radius  $\epsilon_t$  such that the event 280  $\mathcal{E} \coloneqq \{ \| \widehat{\mu}_t - \mu^* \|_{\infty} \le \epsilon_t \ \forall t \in [T] \}$  holds with probability 281 at least  $1 - \delta$ . Second, it defines a parameter  $\beta_t$  such that, 282 conditionally to the realization of the event  $\mathcal{E}$ , the following 283 two conditions hold: (i) the decision space  $\Lambda_{\beta_t}(\hat{\mu}_t)$  contains 284

Algorithm 1 Full-feedback algorithm
function SelectStrategy():
$\boldsymbol{\phi}_t \gets \arg \max_{\boldsymbol{\phi} \in \Lambda_{\beta_t}(\widehat{\boldsymbol{\mu}}_t)} F(\boldsymbol{\phi}, \widehat{\boldsymbol{\mu}}_t)$
return $\phi_t$
function UPDATE( $\boldsymbol{y}_t$ ):
$\widehat{\boldsymbol{\mu}}_{t+1}[\sigma] \leftarrow rac{1}{t} \sum_{ au=1}^t \boldsymbol{y}_{ au}[\sigma] \; orall \sigma \in \Sigma_c$
$\epsilon_{t+1} \leftarrow \sqrt{\frac{\log(2T \Sigma_c /\delta)}{2t}}$
$\beta_{t+1} \leftarrow 2 \mathcal{Z}  \Sigma_r \epsilon_{t+1}$

the optimal signaling scheme  $\phi^*$ ; (ii)  $\Lambda_{2\beta_t}(\mu^*)$  contains the signaling scheme  $\phi_t$ . Intuitively, the first condition is needed to have low sender's regret, while the second one yields signaling schemes that are approximately persuasive.<sup>10</sup>

The polytopal approximation that we provide in Section 4.2 plays a crucial role in the complexity of Algorithm 1. Specifically, it allows it to select the desired  $\phi_t$  in polynomial time by optimizing over the set  $\Lambda_{\beta_t}(\hat{\mu}_t)$ , which can be done efficiently. The use of the set  $\Lambda_{\beta_t}(\hat{\mu}_t)$  over  $\Phi_{\beta_t}^{\diamond}(\hat{\mu}_t)$  is necessary due to the fact that the latter is *not* known to admit an efficient representation. Formally:

**Theorem 4.** Given any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ , Algorithm 1 guarantees:

$$R_T = \mathcal{O}\left(|\mathcal{Z}|\sqrt{T\log(T|\Sigma_c|/\delta)}\right), \quad V_T = \mathcal{O}\left(|\Sigma_r||\mathcal{Z}|\sqrt{T\log(T|\Sigma_c|/\delta)}\right).$$

293 Moreover, the algorithm runs in polynomial time.

### 294 7 Learning with bandit feedback

In this section, we build on Algorithm 1 to deal with bandit feedback, *i.e.*, when at each round  $t \in [T]$ 295 the sender only observes the terminal node  $z_t$  reached at the end of the SDM problem. The main 296 297 difficulties of such a setting can be summarized by the following observations. First, the feedback  $z_t$ only reveals partial information about the prior, and such information also depends on the selected 298 signaling scheme  $\phi_t$ . Second, even if the sender plays a signaling scheme  $\phi \in \Phi$  for an arbitrarily 299 large number of rounds, there is no guarantee that they collect enough information to tell whether 300  $\phi \in \Phi_{\epsilon}^{\diamond}(\mu^{\star})$  or *not* for some  $\epsilon > 0$ . Indeed, the persuasiveness of a signaling scheme depends on 301 all receiver's utilities in the SDM problem, and some parts of the tree may not be reached during a 302

<sup>&</sup>lt;sup>10</sup>See Lemma 9 and 10 in Appendix F for the formal statements of these properties.

sufficiently large number of rounds by committing to  $\phi$ . Thus, any algorithm for the bandit-feedback setting must guarantee a suitable level of exploration over the entire tree, so as to keep track of the entity of the violation of persuasiveness constraints.

We design a two-phase algorithm, whose 306 pseudo-code is provided in Algorithm 2. The 307 algorithm takes as input the number  $N \in [T]$ 308 of rounds devoted to the *first phase* guarantee-309 ing the necessary amount of exploration, as 310 detailed in Section 7.1. During this phase, the 311 SELECTSTRATEGY() procedure implements 312 an efficient deterministic uniform exploration 313 policy, which builds an unbiased estimator 314  $\widehat{\mu}_N$  of  $\mu^*$ . This allows to restrict the space of 315 feasible signaling schemes used in the subse-316 quent phase to those that are approximately 317 persuasive, *i.e.*, those in the set  $\Lambda_{\beta_N}(\widehat{\mu}_N)$ . In 318 Section 7.2, we discuss the second phase of 319 the the algorithm, composed by the rounds 320 t > N, during which the algorithm focuses on 321 the minimization of sender's regret by exploit-322 ing the optimism in face of uncertainty princi-323 324 ple. Finally, in Section 7.3, we provide a lower bound on the trade-off between sender's and 325 receiver's regrets, matching the upper bounds 326 achieved by Algorithm 2 for a large portion 327 of the trade-off frontier. This result formally 328 motivates the necessity of the uniform explo-329

Algorithm 2 Bandit-feedback algorithm

 function SELECTSTRATEGY():

 if 
$$t \leq N$$
 then
 > First Phase

  $\sigma = (h, a) \leftarrow \arg \min_{\sigma \in \Sigma_c} C_t[\sigma]$ 
 > Second Phase

  $\Sigma_s \ni \sigma' \leftarrow \sigma_s(h)$ 
 > Second Phase

 Choose  $\phi_t \in \Phi : \phi_t[\sigma'] = 1$ 
 > Second Phase

  $\phi_t \leftarrow \arg \max_{\phi \in \Lambda_{\beta_N}(\hat{\mu}_N)} \max_{\mu \in C_t(\delta)} F(\phi, \mu)$ 
 return  $\phi_t$ 

 $\begin{aligned} & \text{function UPDATE}(z_t):\\ & \text{Build path } \boldsymbol{p}_t \in \{0,1\}^{|\Sigma_c|} \text{ from } \sigma_c(z_t)\\ & \text{Sample } \boldsymbol{\pi}_t \sim \boldsymbol{\phi}_t \text{ s.t. } \boldsymbol{p}_t[\sigma] = 1 \Rightarrow \sigma \in \Sigma_{\downarrow}(\boldsymbol{\pi}_t)\\ & \text{for } \sigma \in \Sigma_{\downarrow}(\boldsymbol{\pi}_t) \text{ do}\\ & C_{t+1}[\sigma] \leftarrow C_t[\sigma] + 1\\ & \hat{\boldsymbol{\mu}}_{t+1}[\sigma] \leftarrow \frac{1}{C_{t+1}[\sigma]} \sum_{\tau=1}^{C_{t+1}[\sigma]} \boldsymbol{p}_{\tau}[\sigma]\\ & \epsilon_{t+1}[\sigma] \leftarrow \sqrt{\frac{\log(4T|\Sigma_c|/\delta)}{2C_{t+1}[\sigma]}}\\ & \mathcal{C}_{t+1}(\delta) \leftarrow \left\{\boldsymbol{\mu} \left| |\boldsymbol{\mu}[\sigma] - \hat{\boldsymbol{\mu}}_{t+1}[\sigma]| \leq \epsilon_{t+1}[\sigma] \forall \sigma \in \Sigma_c \right. \right\}\\ & \beta_{t+1} \leftarrow 2|\mathcal{Z}||\Sigma_c|\sqrt{\frac{|\Sigma_c|\log(4T|\Sigma_c|/\delta)}{2(t+1)}} \end{aligned}$ 

ration which is performed in the first phase of the algorithm.

#### 331 7.1 Minimizing the receiver's regret

At each round  $t \in [T]$ , the sender observes a terminal node  $z_t \in \mathcal{Z}$  that uniquely determines a path 332 in the tree defining the SDM problem. We encode such a path by means of a vector  $p_t \in \{0,1\}^{|\Sigma_c|}$ 333 such that  $p_t[\sigma] = 1$  if and only if the chance sequence  $\sigma \in \Sigma_c$  lies on the path from the root of 334 the tree to  $z_t$ , namely  $\sigma \preceq \sigma_c(z_t)$ . If the sender commits to a signaling scheme  $\phi_t \in \Phi$ , then it 335 is easy to see that, for every  $\sigma = (h, a) \in \Sigma_c$ , the element  $p_t[\sigma]$  is distributed as a Bernoulli of 336 parameter  $\phi_t[\sigma_s(h)]\mu^*[\sigma]$ . The crucial observation behind the design of our estimator is that, if the 337 sender commits to a deterministic signaling schemes  $\pi_t \in \Pi$  at some round  $t \in [T]$ , then for all the 338 chance sequences  $\sigma \in \Sigma_c$  that are *compatible* with  $\pi_t$ , *i.e.*, that can be observed when  $\pi_t$  is played, 339 we have that  $p_t[\sigma]$  is distributed as a Bernoulli of parameter  $\mu^*[\sigma]$ . Formally, a sequence  $\sigma \in \Sigma_c$ 340 is compatible with  $\pi_t$  if there exists a chance node  $h \in \mathcal{H}_c$  and an outcome  $a \in A(h)$  satisfying 341  $\sigma = (h, a)$  and  $\pi_t[\sigma_s(h)] = 1$ . This observation leads to the following result: 342

**Lemma 3.** For every deterministic signaling scheme  $\pi \in \Pi$ , let

$$\Sigma_{\downarrow}(\boldsymbol{\pi}) \coloneqq \{ \sigma = (h, a) \in \Sigma_c \mid a \in A(h) \land \boldsymbol{\pi}[\sigma_s(h)] = 1 \}$$

Then, during each round  $t \leq N$  of Algorithm 2, it holds  $\mathbb{E}[\mathbf{p}_t[\sigma]] = \boldsymbol{\mu}^*[\sigma]$  for every  $\sigma \in \Sigma_{\downarrow}(\boldsymbol{\pi}_t)$ .

Thus, during the first phase, Algorithm 2 builds the desired estimator  $\hat{\mu}_N$  of  $\mu^*$  as follows. At 345 each round  $t \leq N$ , after observing the feedback  $z_t$ , the algorithm samples a deterministic signaling 346 scheme  $\pi_t \in \Pi$  according to  $\phi_t$  (the one actually selected at t), so that all the sequences  $\sigma \in \Sigma_c$ 347 such that  $p_t[\sigma] = 1$  (or, equivalently,  $\sigma \leq \sigma_c(z_t)$ ) belong to  $\Sigma_{\downarrow}(\pi_t)$ .<sup>11</sup> Then, for every  $\sigma \in \Sigma_{\downarrow}(\pi_t)$ , the algorithm updates the estimator component  $\hat{\mu}_t[\sigma]$  according to  $p_t[\sigma]$ . Since the probability of 348 349 visiting a sequence  $\sigma \in \Sigma_c$  depends on  $\phi_t$  (and, thus, can be arbitrarily small), the first N rounds 350 must be carefully used to ensure that each sequence is explored at least  $N/|\Sigma_c|$  times. To explore 351 a specific sequence  $\sigma \in \Sigma_c$ , we choose a signaling scheme  $\phi_t$  such that  $\sigma \in \Sigma_{\downarrow}(\pi_t)$  for every 352 deterministic  $\pi_t \sim \phi_t$ . The procedure described above is needed for minimizing the receiver's regret, 353 since, in the second phase, the algorithm selects signaling schemes  $\phi_t$  from  $\Lambda_{\beta_N}(\hat{\mu}_N)$ . In particular, 354

<sup>&</sup>lt;sup>11</sup>The sampling can be done efficiently by a straightforward modification of the recursive procedure in [15, 16].

as shown by the following lemma, Algorithm 2 guarantees that the receiver's regret is upper bounded by  $2\beta_N$  at each round t > N, since it defines  $\epsilon_t[\sigma]$  for each sequence  $\sigma \in \Sigma_c$  so that the event  $\tilde{\mathcal{E}} := \{|\boldsymbol{\mu}^*[\sigma] - \hat{\boldsymbol{\mu}}_t[\sigma]| \le \epsilon_t[\sigma] \ \forall (t, \sigma) \in [T] \times \Sigma_c\}$  holds with probability at least  $1 - \delta/2$ .

**Lemma 4.** Under the event  $\tilde{\mathcal{E}}$ , Algorithm 2 guarantees that  $\phi_t \in \Lambda_{2\beta_N}(\mu^*)$  at each round t > N.

#### 359 7.2 Minimizing the sender's regret

Algorithm 2 also needs to guarantee small sender's regret. To do so, we would like that  $\phi^*$  is a valid pick for the algorithm, *i.e.*, it belongs to  $\Lambda_{\beta_N}(\hat{\mu}_t)$ . However, differently from the full-feedback setting, stopping exploration after the first N round does *not* guarantee optimal rates. In order to fix this issue, in the second phase, the algorithm selects  $\phi_t$  optimistically by maximizing the sender's expected utility  $F(\phi, \mu)$  over both  $\phi \in \Lambda_{\beta_N}(\hat{\mu}_N)$  and  $\mu \in C_t(\delta)$ , where  $C_t(\delta)$  is a suitably-defined confidence set centered around  $\hat{\mu}_t$  such that  $\{\mu^* \in C_t(\delta)\} \equiv \tilde{\mathcal{E}}$ , and, thus, it holds with high probability. This guarantees that  $\max_{\mu \in C_t(\delta)} F(\phi^*, \mu) \ge F(\phi^*, \mu^*)$ . Formally:

Lemma 5. If the event  $\tilde{\mathcal{E}}$  holds, then, for every round t > N, it holds that  $\phi^* \in \Lambda_{\beta_N}(\hat{\mu}_t)$  and max $_{\mu \in \mathcal{C}_t(\delta)} F(\phi^*, \mu) \ge F(\phi^*, \mu^*)$ .

Thus,  $F(\phi_t, \mu^*) \approx F(\phi_t, \widehat{\mu}_t) \ge \max_{\mu \in C_t(\delta)} F(\phi^*, \widehat{\mu}) \ge F(\phi^*, \mu^*)$  holds in the limit, implying that  $F(\phi_t, \mu^*)$  converges to  $F(\phi^*, \mu^*)$  after sufficiently many rounds. Formally:

**Theorem 5.** Given any  $\delta \in (0, 1)$  and  $N \in [T]$ , Algorithm 2 guarantees:

$$R_T = \mathcal{O}\left(N + \sqrt{\log(T|\Sigma_c|/\delta)|\Sigma_c|T}\right) \quad and \quad V_T = \mathcal{O}\left(N + T|\mathcal{Z}|\sqrt{|\Sigma_c|\log(T|\Sigma_c|/\delta)/N}\right),$$

with probability at least  $1 - \delta$ . Moreover, the algorithm runs in polynomial time.

In contrast to the case with full feedback, the optimization problem solved by Algorithm 2 belongs to the class of bilinear problems, which are NP-hard in general [18]. However, in Theorem 5 we prove that our specific problem can be solved in polynomial time. Furthermore, notice that Theorem 5 takes as input the number N of rounds devoted to the first phase. Given an  $\alpha \ge 1/2$ , by choosing any  $N = \lfloor T^{\alpha} \rfloor$  we get bounds of  $R_T = \tilde{\mathcal{O}}(T^{\alpha})$  and  $V_T = \tilde{\mathcal{O}}(T^{\max\{\alpha, 1-\frac{\alpha}{2}\}})$ .

#### **7.3 The lower bound frontier**



We conclude by showing that the trade offs between  $V_T$  and  $R_T$  achieved by Algorithm 2 are essentially tight. Previously, we provided an intuition as to why the algorithm needs to uniformly explore the entire tree of the SDM problem. Here, we provide a lower bound that corroborates such a statement. In particular, the following theorem shows that, for any  $\alpha \in [1/2, 1]$ , in order to guarantee a sender's regret of the order of  $\mathcal{O}(T^{\alpha})$ , it is necessary to suffer a receiver's regret of the order of  $\Omega(T^{1-\alpha/2})$ .<sup>12</sup>

**Theorem 6.** For any  $\alpha \in [1/2, 1]$ , there exists a constant  $\gamma \in (0, 1)$  such that no algorithm guarantees both  $R_T = o(T^{\alpha})$  and  $V_T = o(T^{1-\alpha/2})$  with probability greater than  $\gamma$ .

Figure 1 shows on the horizontal axis the order of the T term in  $R_T$ , while, on the vertical axis, 390 it shows the order of the T in  $V_T$ . The shaded area over the blue line shows the achievable trade 391 offs, while the marked red line shows the performances proved in Theorem 5. Thus, we show 392 that Algorithm 2 matches the lower bound for  $\alpha \in [1/2, 2/3]$ . However, when  $\alpha \in [2/3, 1]$ , the 393 guarantees proved in Theorem 5 diverge from the ones proved in the lower bound. This is due to the 394  $N = |T^{\alpha}|$  component in the receiver's regret that becomes dominant when  $\alpha \geq 2/3$ . We conjecture 395 that it is possible to reduce this term to  $\sqrt{N}$ , hence matching the lower bound of Theorem 6. The 396 reason for such a gap between the lower and upper bounds is that, during the first phase, Algorithm 2 397 utilizes signaling schemes without taking into account their persuasiveness, thus incurring in large 398 receiver's regret during the first steps. We leave as future work addressing the question of whether it is 399 possible to design exploration strategies by only using approximately-persuasive signaling schemes. 400

<sup>&</sup>lt;sup>12</sup> For  $\alpha \leq 1/2$ , a simple reduction from a standard multi-armed bandit problem provides a lower bound of  $\Omega(\sqrt{T})$  on both sender's regret  $R_T$  and receiver's regret  $V_T$ .

### 401 **References**

- [1] Ricardo Alonso and Odilon Câmara. Persuading voters. *American Economic Review*, 106(11):
   3590–3605, 2016.
- Itai Arieli and Yakov Babichenko. Private bayesian persuasion. *Journal of Economic Theory*, 182:185–217, 2019.
- [3] Yakov Babichenko and Siddharth Barman. Algorithmic aspects of private Bayesian persuasion.
   In *ITCS*, 2017.
- [4] Umang Bhaskar, Yu Cheng, Young Kun Ko, and Chaitanya Swamy. Hardness results for
   signaling in bayesian zero-sum and network routing games. In *EC*, pages 479–496, 2016.
- [5] Peter Bro Miltersen and Or Sheffet. Send mixed signals: earn more, work less. In *EC*, pages
   234–247, 2012.
- [6] Ozan Candogan. Persuasion in networks: Public signals and k-cores. In *EC*, pages 133–134, 2019.
- [7] Matteo Castiglioni and Nicola Gatti. Persuading voters in district-based elections. In AAAI,
   pages 5244–5251, 2021.
- [8] Matteo Castiglioni, Andrea Celli, and Nicola Gatti. Persuading voters: It's easy to whisper, it's
   hard to speak loud. In AAAI, pages 1870–1877, 2020.
- [9] Matteo Castiglioni, Andrea Celli, Alberto Marchesi, and Nicola Gatti. Online Bayesian
   persuasion. In *NeurIPS*, pages 16188–16198, 2020.
- [10] Matteo Castiglioni, Andrea Celli, Alberto Marchesi, and Nicola Gatti. Signaling in Bayesian
   network congestion games: the subtle power of symmetry. In AAAI, pages 5252–5259, 2021.
- [11] Matteo Castiglioni, Alberto Marchesi, Andrea Celli, and Nicola Gatti. Multi-receiver online
   bayesian persuasion. In *ICML*, volume 139, pages 1314–1323, 2021.
- [12] Andrea Celli, Stefano Coniglio, and Nicola Gatti. Private bayesian persuasion with sequential
   games. AAAI, 34(02):1886–1893, 2020.
- [13] Andrea Celli, Alberto Marchesi, Gabriele Farina, and Nicola Gatti. No-regret learning dynamics
   for extensive-form correlated equilibrium. *NeurIPS*, 33:7722–7732, 2020.
- [14] Nicolo Cesa-Bianchi and Gábor Lugosi. *Prediction, learning, and games*. Cambridge university
   press, 2006.
- [15] Gabriele Farina, Andrea Celli, Alberto Marchesi, and Nicola Gatti. Simple uncoupled no-regret
   learning dynamics for extensive-form correlated equilibrium. *arXiv preprint arXiv:2104.01520*,
   2021.
- [16] Gabriele Farina, Robin Schmucker, and Tuomas Sandholm. Bandit linear optimization for
   sequential decision making and extensive-form games. In *AAAI*, volume 35, pages 5372–5380,
   2021.
- [17] J Gan, R Majumdar, G Radanovic, and A Singla. Bayesian persuasion in sequential decision making. In AAAI, 2022.
- [18] Christopher J Hillar and Lek-Heng Lim. Most tensor problems are np-hard. *Journal of the* ACM, 60(6):1–39, 2013.
- [19] Emir Kamenica and Matthew Gentzkow. Bayesian persuasion. *American Economic Review*,
   101(6):2590–2615, 2011.
- [20] Daphne Koller, Nimrod Megiddo, and Bernhard Von Stengel. Efficient computation of equilibria
   for extensive two-person games. *Games and economic behavior*, 14(2):247–259, 1996.
- [21] Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- [22] Dustin Morrill, Ryan D'Orazio, Reca Sarfati, Marc Lanctot, James R Wright, Amy R Greenwald,
   and Michael Bowling. Hindsight and sequential rationality of correlated play. In *AAAI*,
   volume 35, pages 5584–5594, 2021.
- [23] Zinovi Rabinovich, Albert Xin Jiang, Manish Jain, and Haifeng Xu. Information disclosure as a
   means to security. In AAMAS, pages 645–653, 2015.
- [24] Bernhard Von Stengel and Françoise Forges. Extensive-form correlated equilibrium: Definition
   and computational complexity. *Mathematics of Operations Research*, 33(4):1002–1022, 2008.

- Jibang Wu, Zixuan Zhang, Zhe Feng, Zhaoran Wang, Zhuoran Yang, Michael I. Jordan, and
   Haifeng Xu. Sequential information design: Markov persuasion process and its efficient
   reinforcement learning. *arXiv preprint arXiv:2202.10678*, 2022.
- [26] Haifeng Xu, Rupert Freeman, Vincent Conitzer, Shaddin Dughmi, and Milind Tambe. Signaling
   in Bayesian Stackelberg games. In *AAMAS*, pages 150–158, 2016.
- [27] You Zu, Krishnamurthy Iyer, and Haifeng Xu. Learning to persuade on the fly: Robustness against ignorance. In *EC*, pages 927–928, 2021.

## 459 Checklist

460	1. For all authors
461 462	<ul> <li>(a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]</li> </ul>
463 464	(b) Did you describe the limitations of your work? [Yes] The setting is based on precise assumptions. If those assumptions are not met then our techiques are not applicable.
465 466	(c) Did you discuss any potential negative societal impacts of your work? [N/A] Our work is mailny theoretical and thus the question does not apply.
467 468	<ul><li>(d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]</li></ul>
469	2. If you are including theoretical results
470	(a) Did you state the full set of assumptions of all theoretical results? [Yes]
471 472 473	(b) Did you include complete proofs of all theoretical results? [Yes] While all the main statements are contained in the main paper, all the proof are deferred to the appendix due to space constraints.
474	3. If you ran experiments
475 476	<ul> <li>(a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [N/A]</li> </ul>
477 478	<ul><li>(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [N/A]</li></ul>
479 480	(c) Did you report error bars (e.g., with respect to the random seed after running experi- ments multiple times)? [N/A]
481 482	<ul><li>(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [N/A]</li></ul>
483	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
484	(a) If your work uses existing assets, did you cite the creators? [N/A]
485	(b) Did you mention the license of the assets? [N/A]
486 487	(c) Did you include any new assets either in the supplemental material or as a URL? $[N/A]$
488 489	(d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
490 491	(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
492	5. If you used crowdsourcing or conducted research with human subjects
493 494	<ul> <li>(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]</li> </ul>
495 496	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
497 498	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

### **A Relation with the classical Bayesian persuasion framework**

In this section, we clarify the relationship be-tween the BPSDM problems that we study

<sup>502</sup> in this paper and the classical Bayesian per-

suasion framework introduced by Kamenicaand Gentzkow [19]. In particular, we show

that any instance of the classical Bayesian

- <sup>506</sup> persuasion problem can be mapped to an
- instance of the BPSDM problem.

508 A Bayesian persuasion problem instance is

- 509 defined by a set  $\mathcal{A}$  of  $k \coloneqq |\mathcal{A}|$  actions for
- the receiver, a set S of signals for the sender,
- and a set  $\Theta$  of  $d \coloneqq |\Theta|$  possible outcomes



Figure 2: Instance of BPSDM problem corresponding to a given instance of Bayesian persuasion problem.

of a (single) random event (called *states of nature* in the classical Bayesian persuasion terminology).

The receiver's payoff function is  $u^{\mathsf{R}} : \Theta \times \mathcal{A} \to [0, 1]$ , while the sender's one is  $u^{\mathsf{S}} : \Theta \times \mathcal{A} \to [0, 1]$ . The sender observes the realized state of nature, which is drawn according to a commonly-known prior distribution  $\mu \in \Delta_{\Theta}$ . Then, they partially disclose information about the state by committing to

a signaling scheme  $\varphi: \Theta \to \Delta_S$ , which is a randomized mapping from states of nature to signals for the receiver. Thus, the interaction between the sender and the receiver is as follows:

(i) The sender commits to a publicly known signaling scheme  $\varphi$ .

(ii) The sender observes the state of nature  $\theta \sim \mu$ .

(iii) The sender samples a signal  $s \sim \varphi(\theta, \cdot)$  and sends it to the receiver.

(iv) The receiver computes their posterior belief over the states  $\Theta$ .

(v) The receiver plays an action  $a \in A$  that maximizes their expected payoff.

The posterior beliefs that the receiver computes in step (iv) after observing a signal  $s \in S$  are defined by a probability distribution  $\xi_s \in \Delta_{\Theta}$  such that:

$$\xi_s( heta) \coloneqq rac{\mu( heta) arphi( heta,s)}{\sum_{ heta' \in \Theta} \mu( heta') arphi( heta',s)},$$

and, thus, after observing signal s the receiver plays an action

$$a \in \arg \max_{a' \in \mathcal{A}} \sum_{\theta \in \Theta} \xi_s(\theta) u^{\mathsf{R}}(\theta, a')$$

A revelation-principle-style argument [19] allow the sender to focus on direct and persuasive signaling schemes, where the latter property means that  $S \equiv A$ , with signals corresponding to actions recommendations for the receiver. A persuasive signaling scheme  $\varphi : \Theta \rightarrow \Delta_S$  is such that the receiver is always incentivized to follow action recommendations; formally:

$$\sum_{\theta \in \Theta} \mu(\theta)\varphi(\theta, a)u^{\mathbb{R}}(\theta, a) \ge \sum_{\theta \in \Theta} \mu(\theta)\varphi(\theta, a)u^{\mathbb{R}}(\theta, a') \quad \forall a, a' \in \mathcal{A}.$$
 (2)

Instance mapping. Given an instance of the classical Bayesian persuasion problem [19], a corresponding (equivalent) instance of our BPSDM problem can be constructed as follows:

- (1) There is a unique chance node  $h_0$  which is the root of the tree defining the SDM problem.
- (2) At the chance node, there are *d* possible outcomes (namely  $A(h_0) \equiv \Theta$ ), each corresponding to a state of nature  $\theta \in \Theta$  and having probability  $\mu(\theta)$  of occurring, so that with a slight abuse of notation we can write  $\mu^*[\varnothing] = 1$  and  $\mu^*[\theta] = \mu(\theta)$  for all  $\theta \in \Theta$ .
- (3) The receiver has a unique infoset *I*, which contains one decision node for each possible
   outcome at the chance node.
- (4) At infoset I, the receiver has a set  $A(I) \equiv A$  of available actions.
- (5) Terminal nodes  $\mathcal{Z}$  are determined by state of nature, receiver's action pairs, so that each  $\theta_i \in \Theta$  and  $a_j \in \mathcal{A}$  define a corresponding terminal node  $z_{i,j}$  in the SDM problem.

The following theorem formally states that our definition of persuasiveness (Definition 2) instantiated to the BPSDM problem instances described above is equivalent to the definition of persuasiveness for classical Bayesian persuasion problems (Equation (2)). This establishes that our framework

encompasses classical Bayesian persuasion problems as a special case.

Theorem 7. Given any Bayesian persuasion instance, a signaling scheme is persuasive (Equation (2))
 if and only if it is persuasive (Definition 2) in the corresponding instance of BPSDM problem.

*Proof.* It is sufficient to prove the equivalence between Equation (1) for  $\epsilon = 0$  and Equation (2) applied to the BPSDM problem instance representing the given Bayesian persuasion instance. To do that, we employ Theorem 1 and Lemma 6 in such a BPSDM problem instance, so that, using the notation introduced in this section, it is straightforward to see that Equation (1) reads as follows:

$$\max_{a' \in \mathcal{A}} \sum_{\theta \in \Theta} \varphi(\theta, a) \mu(\theta) u^{\mathsf{R}}(\theta, a') - \sum_{\theta \in \Theta} \varphi(\theta, a) \mu(\theta) u^{\mathsf{R}}(\theta, a) \le 0 \quad \forall a \in \mathcal{A}$$

<sup>551</sup> By rearranging the terms, we get Equation (2), which concludes the proof.

### 552 **B** Example of SDM problem and its sets of sequences

Figure 3 shows a simple instance of a SDM problem. This is defined by a tree whose set of chance nodes is  $\mathcal{H}_c = \{h_0\}$ , while the set of decision nodes is  $\mathcal{H}_d = \{h_1, h_2, h_3\}$ . The set of terminal nodes is  $\mathcal{Z} = \{z_1, \dots, z_6\}$ . Moreover, the set of decision nodes  $\mathcal{H}_d$  is partitioned into the set partition  $\mathcal{I} = \{I, J\}$ , which is made by two infosets  $I = \{h_1\}$  and  $J = \{h_2, h_3\}$ .



Figure 3: Example of SDM problem and its sets of sequences  $\Sigma_r$ ,  $\Sigma_s$ , and  $\Sigma_c$ .

The sets of sequences are constructed as follows. For the chance agent, we have that  $\Sigma_c = \{(h_0, a), (h_0, b), (h_0, c)\}$ , while for the receiver we have that  $\Sigma_r = \{(I, d), (I, e), (J, f), (J, g)\}$ . Let us remark that, since the receiver cannot distinguish between nodes  $h_2$  and  $h_3$ , it only has 2 sequences originating from such nodes; namely (J, f) and (J, g). Finally, the sender can be thought of as a perfect-information agent selecting action recommendations for the receiver at decision nodes, so that their set of sequences is  $\Sigma_s = \{(h_1, d), (h_1, e), (h_2, f), (h_2, g), (h_3, f), (h_3, g)\}$ .

### 563 C Additional notation needed in the proofs

In this section, we introduce some additional notation that will be useful in the proofs.

We denote by  $\Pi_r := \mathcal{X}_r \cap \{0,1\}^{|\Sigma_r|}$  the set of *deterministic* sequence-form strategies (a.k.a. *pure strategies*) of the receiver, which are the strategies specifying to play a single action with probability one at each infoset. The set of receiver's deterministic strategies in the SDM sub-problem that starts from an infoset  $I \in \mathcal{I}$  is denoted as  $\Pi_{r,I} := \mathcal{X}_{r,I} \cap \{0,1\}^{|\Sigma_r|}$ . Moreover, we let  $\Sigma_{r,I} \subseteq \Sigma_r$  be the set of receiver's sequences in the SDM sub-problem that starts from an infoset  $I \in \mathcal{I}$ ; formally,  $\Sigma_{r,I} := \{\sigma \in \Sigma_r \mid \sigma_r(I) \preceq \sigma \land \exists z \in \mathcal{Z}(I) : \sigma \preceq \sigma_r(z)\}$ 

Given any infoset  $I \in \mathcal{I}$ , we let  $\mathcal{Z}(I) \subset \mathcal{Z}$  be the set of terminal nodes  $z \in \mathcal{Z}$  such that the path from the root of the tree to z passes through a node in I. Moreover, given  $\sigma = (I, a)$  with  $a \in A(I)$ , we define  $\mathcal{Z}(\sigma) = \mathcal{Z}(I, a) \subset \mathcal{Z}(I)$  as the set of terminal nodes whose corresponding paths include playing action a at a node in I. For every infoset  $I \in \mathcal{I}$ , we also introduce a function  $h_I : Z(I) \to I$ such that  $h_I(z)$  defines the unique node  $h \in I$  on the path from the root of the tree to z. Given an infoset  $I \in \mathcal{I}$  and an action  $a \in A(I)$ , we define  $C(I, a) \subseteq \mathcal{I}$  as the set of all the infosets which immediately follow infoset I through action a, *i.e.*, those infosets  $J \in \mathcal{I}$  such that  $\sigma_r(J) = (I, a)$ . Moreover, we let  $C(I) \subseteq \mathcal{I}$  be the set of all infosets that follow  $\mathcal{I}$ , *i.e.*, those infosets  $J \in \mathcal{I}$  such that there exits  $a \in A(I)$  with  $\sigma = (I, a)$  such that  $\sigma \preceq \sigma_r(J)$ .

### 580 **D Proofs omitted from Section 4**

Let us remark that all the results in Section 4 can be straightforwardly generalized to the case of a generic  $\mu \in \mathcal{X}_c$ , as needed for the proofs of the results in Sections 6 and 7. For ease of exposition, we state and prove the results of Section 4 for the prior  $\mu^*$ .

First, we prove a preliminary lemma that allows us to express the receiver's expected utility difference between using a  $(\sigma, \rho_{\sigma})$ -SPDP and following action recommendations by only considering the terminal nodes under the infoset in which the SPDP prescribed to deviate. A similar result for the case of correlated strategies can be found in [13, Appendix A].

**Lemma 6.** Given  $\phi \in \Phi$ , for every  $(\sigma, \rho_{\sigma})$ -SPDP with  $\sigma = (I, a) \in \Sigma_r$  and  $\rho_{\sigma} \in \mathcal{X}_{r,I}$ , it holds:

$$U_{\sigma \to \boldsymbol{\rho}_{\sigma}}(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) = \sum_{z \in Z(I)} \boldsymbol{\phi}[(h_{I}(z), a)] \boldsymbol{\rho}_{\sigma}[\sigma_{r}(z)] \boldsymbol{\mu}^{\star}[\sigma_{c}(z)] u(z) + -\sum_{z \in Z(\sigma)} \boldsymbol{\phi}[\sigma_{s}(z)] \boldsymbol{\mu}^{\star}[\sigma_{c}(z)] u(z).$$

*Proof.* We define the following three disjoint events for any  $(\sigma, \rho_{\sigma})$ -SPDP, where  $\sigma = (I, a)$ .

590 (C1): A terminal node  $z \in \mathcal{Z}(\sigma)$  is reached.

(C2): A terminal node  $z \in \mathcal{Z}(I, a')$  for some  $a' \neq a \in A(I)$  is reached.

592 (C3): A terminal node  $z \in \mathcal{Z}/\mathcal{Z}(I)$  is reached.

l

Next, under each event, we define the probability  $p_{\sigma \to \rho_{\sigma}}(z)$  of reaching a terminal node z:

(C1): Since  $z \in \mathcal{Z}(\sigma)$ , the node z is reached by means of the continuation strategy  $\rho_{\sigma}$ . Thus:  $p_{\sigma \to \rho}^{(1)}(z) \coloneqq \phi[(h_I(z), a)] \mu^*[\sigma_c(z)] \rho_{\sigma}[\sigma_r(z)].$ 

(C2): Since  $z \in \mathcal{Z}(I, a')$  for  $a' \neq a \in A(I)$ , the node z can be reached either by deviating and then committing to the continuation strategy  $\rho_{\sigma}$  or by following recommendations. Moreover, these two cases are exclusive, and, thus, we can write:

$$p_{\sigma \to \rho_{\sigma}}^{(2)}(z) \coloneqq \phi[(h_I(z), a)] \boldsymbol{\mu}^*[\sigma_c(z)] \boldsymbol{\rho}_{\sigma}[\sigma_r(z)] + \phi[\sigma_s(z)] \boldsymbol{\mu}^*[\sigma_c(z)].$$

(C3): Since 
$$z \in \mathbb{Z}/\mathbb{Z}(I)$$
, the node z is reached by following recommendations:

$$p_{\sigma \to \rho_{\sigma}}^{(3)}(z) \coloneqq \boldsymbol{\phi}[\sigma_s(z)] \boldsymbol{\mu}^{\star}[\sigma_c(z)].$$

We observe that  $p_{\sigma \to \rho_{\sigma}}^{(2)}(z) = p_{\sigma \to \rho_{\sigma}}^{(1)}(z) + p_{\sigma \to \rho_{\sigma}}^{(3)}(z)$ , and, thus, we can write  $U_{\sigma \to \rho_{\sigma}}(\phi, \mu^{\star})$  as:  $U_{\sigma \to \rho_{\sigma}}(\phi, \mu^{\star}) \coloneqq \sum_{z \in \mathcal{Z}(\sigma)} p_{\sigma \to \rho_{\sigma}}^{(1)}(z)u(z) + \sum_{\substack{z \in \mathcal{Z}(I,a'):\\a' \neq a \in A(I)}} p_{\sigma \to \rho_{\sigma}}^{(2)}(z)u(z) + \sum_{z \in \mathcal{Z}/\mathcal{Z}(\sigma)} p_{\sigma \to \rho_{\sigma}}^{(3)}(z)u(z)$ .

Furthermore, by using the definition of  $p_{\sigma \to \rho_{\sigma}}^{(3)}(z)$ , we can write  $U(\phi, \mu) \coloneqq \sum_{z \in \mathcal{Z}} u(z) p_{\sigma \to \rho_{\sigma}}^{(3)}(z)$ . Thus:

$$U_{\sigma \to \boldsymbol{\rho}_{\sigma}}(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) = \sum_{z \in \mathcal{Z}(I)} p_{\sigma \to \boldsymbol{\rho}_{\sigma}}^{(1)}(z)u(z) - \sum_{z \in \mathcal{Z}(\sigma)} p_{\sigma \to \boldsymbol{\rho}_{\sigma}}^{(3)}(z)u(z)$$

which is the statement of the lemma by substituting the definitions of  $p_{\sigma \to \rho_{\sigma}}^{(1)}(z)$  and  $p_{\sigma \to \rho_{\sigma}}^{(3)}(z)$ .

- Now, we exploit Lemma 6 to prove the following local decomposition of a DP into SPDPs. 603
- **Theorem 1.** Given a signaling scheme  $\phi \in \Phi$  and a  $(\omega, \rho)$ -DP, it holds: 604

$$U^{\boldsymbol{\omega} \to \boldsymbol{\rho}}(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) = \sum_{\sigma \in \Sigma_r} \boldsymbol{\omega}[\sigma] \Big( U_{\sigma \to \boldsymbol{\rho}_{\sigma}}(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) \Big).$$

*Proof.* For any terminal node  $z \in \mathbb{Z}$ , let  $p^{\omega \to \rho}(z; \phi, \mu^{\star})$  be the probability of reaching node z when 605 the receiver employs the  $(\omega, \rho)$ -DP under the signaling scheme  $\phi$  and the prior  $\mu^*$ . It holds: 606

$$p^{\boldsymbol{\omega} \to \boldsymbol{\rho}}(z; \boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) \coloneqq \sum_{\sigma = (I, a) \in \Sigma_r : \sigma \preceq \sigma_r(z)} \boldsymbol{\omega}[\sigma] \boldsymbol{\phi}[(h_I(z), a)] \boldsymbol{\rho}_{\sigma}[\sigma_r(z)] \boldsymbol{\mu}^{\star}[\sigma_c(z)] + \phi[\sigma_s(z)] \boldsymbol{\mu}^{\star}[\sigma_c(z)] \left(1 - \sum_{\sigma \in \Sigma_r : \sigma \preceq \sigma_r(z)} \boldsymbol{\omega}[\sigma]\right).$$

The sum in the first term in the definition of  $p^{\omega \to \rho}(z; \phi, \mu^{\star})$  accounts for the probabilities of reaching 607 z when the receiver reaches infoset I, is recommended to play action a, and deviates by following the 608 continuation strategy  $\rho_{\sigma}$  thereafter, for all the sequences  $\sigma = (I, a)$  that precede the sequence  $\sigma_r(z)$  reaching z. Instead, the second term in the definition of  $p^{\omega \to \rho}(z; \phi, \mu^*)$  accounts for the probability of reaching z by following recommendations. Thus,  $U^{\omega \to \rho}(\phi, \mu^*) = \sum_{z \in \mathbb{Z}} p^{\omega \to \rho}(z; \phi, \mu^*)u(z)$ . 609 610 611

By rearranging the terms in  $U^{\omega \to \rho}(\phi, \mu^{\star})$ , we get to the following result: 612

$$U^{\boldsymbol{\omega} \to \boldsymbol{\rho}}(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) = U(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) + \sum_{z \in \mathcal{Z}} \left[ \sum_{\sigma = (I, a): \sigma \preceq \sigma_r(z)} \boldsymbol{\omega}[\sigma] \boldsymbol{\phi}[(h_I(z), a)] \boldsymbol{\rho}_{\sigma}[\sigma_r(z)] \boldsymbol{\mu}^{\star}[\sigma_c(z)] u(z) + - \sum_{\sigma \in \Sigma_r: \sigma \preceq \sigma_r(z)} \boldsymbol{\omega}[\sigma] \boldsymbol{\phi}[\sigma_s(z)] \boldsymbol{\mu}^{\star}[\sigma_c(z)] u(z) \right] \right]$$
$$= U(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) - \sum_{\sigma \in \Sigma_r} \boldsymbol{\omega}[\sigma] \sum_{z \in \mathcal{Z}(\sigma)} \boldsymbol{\phi}[\sigma_s(z)] \boldsymbol{\mu}^{\star}[\sigma_c(z)] u(z) + + \sum_{\sigma \in \Sigma_r} \boldsymbol{\omega}[\sigma] \sum_{z \in \mathcal{Z}(I)} \boldsymbol{\phi}[(h_I(z), a)] \boldsymbol{\rho}_{\sigma}[\sigma_r(z)] \boldsymbol{\mu}^{\star}[\sigma_c(z)] u(z).$$
(3)

Thus, by combining Lemma 6 with Equation (3) we get that: 613

$$U^{\boldsymbol{\omega} \to \boldsymbol{\rho}}(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) = \sum_{\sigma \in \Sigma_{r}} \boldsymbol{\omega}[\sigma] \left[ U_{\sigma \to \boldsymbol{\rho}_{\sigma}}(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) \right],$$
  
udes the proof.

- which concludes the proof. 614
- **Corollary 1.** *Given a signaling scheme*  $\phi \in \Phi$ *, the following holds:* 615

$$\max_{(\boldsymbol{\omega},\boldsymbol{\rho})\in\Omega\times\mathcal{P}} U^{\boldsymbol{\omega}\to\boldsymbol{\rho}}(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) \leq \sum_{\sigma=(I,a)\in\Sigma_r} \left[ \max_{\boldsymbol{\rho}_{\sigma}\in\mathcal{X}_{r,I}} U_{\sigma\to\boldsymbol{\rho}_{\sigma}}(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) \right]^{+}.$$

*Proof.* By using Theorem 1, we derive the following: 616

$$\begin{aligned} \max_{(\boldsymbol{\omega},\boldsymbol{\rho})\in\Omega\times\mathcal{P}} U^{\boldsymbol{\omega}\to\boldsymbol{\rho}}(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) &= \max_{(\boldsymbol{\omega},\boldsymbol{\rho})\in\Omega\times\mathcal{P}}\sum_{\sigma\in\Sigma_{r}} \boldsymbol{\omega}[\sigma] \left( U_{\sigma\to\boldsymbol{\rho}_{\sigma}}(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) \right) \\ &\leq \max_{(\boldsymbol{\omega},\boldsymbol{\rho})\in\Omega\times\mathcal{P}}\sum_{\sigma\in\Sigma_{r}} \boldsymbol{\omega}[\sigma] \left[ U_{\sigma\to\boldsymbol{\rho}_{\sigma}}(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) \right]^{+} \\ &\leq \max_{\boldsymbol{\rho}\in\mathcal{P}}\sum_{\sigma\in\Sigma_{r}} \left[ U_{\sigma\to\boldsymbol{\rho}_{\sigma}}(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) \right]^{+} \\ &= \sum_{\sigma\in\Sigma_{r}} \left[ \max_{\boldsymbol{\rho}_{\sigma}\in\mathcal{X}_{r,I}} U_{\sigma\to\boldsymbol{\rho}_{\sigma}}(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) \right]^{+}. \end{aligned}$$

This concludes the proof. 617

#### **Lemma 1.** The set $\Lambda_{\epsilon}(\mu^{\star})$ can be described by means of a polynomial number of linear constraints.

- Proof. In order to prove that the set  $\Lambda_{\epsilon}(\mu^{\star})$  can be described by means of linear constraints, we employ duality arguments related to the max problem in the definition of  $\Lambda_{\epsilon}(\mu^{\star})$  (Definition 4).
- By Lemma 6, for every sequence  $\sigma = (I, a) \in \Sigma_r$ , we can rewrite the expression in the left-hand side of the inequality characterizing  $\Lambda_{\epsilon}(\mu^{\star})$  in Definition 4 as follows:

$$\max_{\boldsymbol{\rho}_{\sigma}\in\mathcal{X}_{r,I}}\left\{\sum_{z\in Z(I)}\phi[(h_{I}(z),a)]\boldsymbol{\rho}_{\sigma}[\sigma_{r}(z)]\boldsymbol{\mu}^{\star}[\sigma_{c}(z)]u(z)\right\}-\sum_{z\in Z(\sigma)}\phi[\sigma_{s}(z)]\boldsymbol{\mu}^{\star}[\sigma_{c}(z)]u(z),$$

so that  $\Lambda_{\epsilon}(\mu^{\star})$  can be expressed as the set of all  $\phi \in \Phi$  such that the above expression has value less than or equal to  $\epsilon/|\Sigma_r|$  for every  $\sigma \in \Sigma_r$ . Observe that the expression in the max operator is a linear function of  $\rho_{\sigma}$ , and that the set  $\mathcal{X}_{r,I}$  is a polytope by definition. Thus, for every  $\sigma = (I, a) \in \Sigma_r$ , the maximization above can be equivalently rewritten as the following linear program:

$$\max_{\boldsymbol{x}^{I,a} \ge \boldsymbol{0}} \left( \boldsymbol{x}^{I,a} \right)^{\top} \boldsymbol{c}(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) \quad \text{s.t.}$$
(4a)

$$\boldsymbol{F}_{I}\boldsymbol{x}^{I,a} = \boldsymbol{f}_{I} \tag{4b}$$

where  $\boldsymbol{x}^{I,a}$  is a vector of variables indexed over sequences  $\Sigma_{r,I} \cup \{\sigma_r(I)\}$ . Notice that  $c(\phi, \boldsymbol{\mu}^*) \in \mathbb{R}^{|\Sigma_{r,I}|}$  is a vector of coefficients such that the component corresponding to each  $\sigma' \in \Sigma_{r,I}$  is

$$c(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star})[\boldsymbol{\sigma}'] \coloneqq \sum_{z \in \mathcal{Z}(I): \sigma_r(z) = \sigma'} \boldsymbol{\phi}[(h_I(z), a)] \boldsymbol{\mu}^{\star}[\sigma_c(z)] u(z),$$

while  $c(\phi, \mu^{\star})[\sigma_r(I)] := 0$ . Moreover,  $F_I \in \{-1, 0, 1\}^{(1+|C(I)|) \times |\Sigma_{r,I}|}$  is a matrix of coefficients whose components are defined as follows:  $[F_I]_{I_{\varnothing},\sigma_r(I)} := 1$  and  $[F_I]_{I_{\varnothing},\sigma'} := 0$  for all sequences  $\sigma' \in \Sigma_{r,I}$ , where  $I_{\varnothing}$  is a fictitious infoset indexing the first row, while, for every infoset  $J \in C(I)$ following I (this included) and sequence  $\sigma' \in \Sigma_{r,I} \cup \{\sigma_r(I)\}$ :

$$[\mathbf{F}_I]_{J,\sigma'} \coloneqq \begin{cases} -1 & \text{if} \quad \sigma' = \sigma_r(J) \\ 1 & \text{if} \quad \sigma' = (J,a') \text{ for some } a' \in A(J) \\ 0 & \text{ otherwise} \end{cases}$$

- Finally,  $f_I \in \{0,1\}^{1+|C(I)|}$  is a vector whose components are all zero apart from that one corresponding to the sequence  $\sigma_r(I)$ , which is one (see also [20]).
- <sup>635</sup> The dual linear program of Problem (4) reads as:

$$\min_{\boldsymbol{y}^{I,a}} \boldsymbol{y}^{I,a}[I_{\varnothing}] \quad \text{s.t.}$$
(5a)

$$\boldsymbol{F}_{I}^{\top}\boldsymbol{y}^{I,a} \ge \boldsymbol{c}(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}), \tag{5b}$$

where  $y^{I,a}$  is a vector of dual variables indexed over  $C(I) \cup \{I_{\emptyset}\}$ . For ease of notation, we let OPT<sub>I,a</sub> be the optimal value of Problem (5) instantiated for the sequence  $\sigma = (I, a)$ .

<sup>638</sup> By strong duality, we have that the optimal value of the primal (Problem (4)) is equal to the optimal <sup>639</sup> value of the dual (Problem (5)), and this allows us to readily rewrite the set  $\Lambda_{\epsilon}(\mu^{\star})$  as follows:

$$\Lambda_{\epsilon}(\boldsymbol{\mu}^{\star}) = \left\{ \boldsymbol{\phi} \in \Phi \ \Big| \ \mathsf{OPT}_{I,a} - \sum_{z \in Z(\sigma)} \boldsymbol{\phi}[\sigma_s(z)] \boldsymbol{\mu}^{\star}[\sigma_c(z)] u(z) \le \frac{\epsilon}{|\Sigma_r|} \quad \forall \sigma = (I,a) \in \Sigma_r \right\}.$$
(6)

Moreover, we can remove  $OPT_{I,a}$  in Equation (6) since it appears in in the right-hand side of a inequality and Problem (5) is a min problem. Thus, the set  $\Lambda_{\epsilon}(\mu^{\star})$  can be written as follows:

$$\Lambda_{\epsilon}(\boldsymbol{\mu}^{\star}) = \left\{ \boldsymbol{\phi} \in \Phi \ \Big| \ \exists \boldsymbol{y}^{I,a} \in \mathbb{R}^{1+|C(I)|} : \boldsymbol{y}^{I,a}[I_{\varnothing}] - \sum_{z \in Z(\sigma)} \boldsymbol{\phi}[\sigma_{s}(z)] \boldsymbol{\mu}^{\star}[\sigma_{c}(z)] u(z) \leq \frac{\epsilon}{|\Sigma_{r}|} \\ \wedge \boldsymbol{F}_{I}^{\top} \boldsymbol{y}^{I,a} \geq \boldsymbol{c}(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) \quad \forall \sigma = (I,a) \in \Sigma_{r} \right\},$$

$$(7)$$

which is comprised of a polynomial number of inequalities and variables, concluding the proof.

Let us also notice that, by expanding the constraints of Problem (5), one can easily check that they can be equivalently rewritten recursively, as follows. For every sequence  $\sigma' = (J, a') \in \Sigma_{r,I}$ , Constraints (5b) can be rewritten as:

$$\boldsymbol{y}^{I,a}[J] \ge \sum_{z \in \mathcal{Z}(I):\sigma_r(z)=(J,a')} \boldsymbol{\phi}[(h_I(z),a)] \boldsymbol{\mu}^{\star}[\sigma_c(z)] u(z) + \sum_{K \in C(J,a')} \boldsymbol{y}^{I,a}[K],$$
(8)

while, for sequence  $\sigma_r(I)$ , Constraint (5b) can be written as  $y^{I,a}[I_{\varnothing}] \ge y^{I,a}[I]$ . Intuitively, at any optimal solution to Problem (5), we can interpret the value of the dual variable  $y^{I,a}[I_{\varnothing}]$  as the receiver's expected utility obtained by playing the best possible continuation strategy after being recommended action a at infoset I. Indeed, the first term in the right-hand-side of Equation (8) is the utility immediately obtainable after playing a' at infoset J, while the second term recursively encodes the utilities obtained (non-immediately) following a' at J.

- **Lemma 2.** It is always the case that  $\Phi^{\diamond}(\mu^{\star}) \equiv \Lambda(\mu^{\star}) \subseteq \Lambda_{\epsilon}(\mu^{\star}) \subseteq \Phi_{\epsilon}^{\diamond}(\mu^{\star})$ .
- *Proof.* First, we prove that  $\Phi^{\diamond}(\mu^{\star}) \equiv \Lambda(\mu^{\star})$ . Suppose that  $\phi \in \Phi^{\diamond}(\mu^{\star})$ . Then, Definition 2 implies

$$U^{\boldsymbol{\omega} \to \boldsymbol{\rho}}(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) \leq 0,$$

for every  $\omega \in \Omega$  and  $\rho \in \mathcal{P}$ . Thus, by Theorem 1 we have that:

$$\sum_{\sigma \in \Sigma_r} \boldsymbol{\omega}[\sigma] \Big( U_{\sigma \to \boldsymbol{\rho}_{\sigma}}(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) \Big) \leq 0,$$

for every  $\omega \in \Omega$  and  $\rho \in \mathcal{P}$ , which implies that:

$$\max_{\boldsymbol{\rho}_{\sigma}\in\mathcal{X}_{r,I}} U_{\sigma\to\boldsymbol{\rho}_{\sigma}}(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) \leq 0 \quad \forall \sigma \in \Sigma_{r},$$

- and  $\phi \in \Lambda(\mu^*)$ , proving the first part of the statement.
- <sup>657</sup> On the other hand,  $\Lambda(\mu^*) \subseteq \Phi^{\diamond}(\mu^*)$  is directly implied by Corollary 1. Thus,  $\Lambda(\mu^*) \equiv \Phi^{\diamond}(\mu^*)$ . <sup>658</sup> Moreover, from Definition 4 it trivially follows that  $\Lambda(\mu^*) \subseteq \Lambda_{\epsilon}(\mu^*)$ .
- $\operatorname{Horeover}, \operatorname{Hom} \operatorname{Definition} + \operatorname{Retriving follows that } \operatorname{H}(\mu^{-}) \subseteq \operatorname{H}_{\mathcal{E}}(\mu^{-}).$
- Finally, we prove that  $\Lambda_{\epsilon}(\mu^{\star}) \subseteq \Phi_{\epsilon}^{\diamond}(\mu^{\star})$ . Given  $\epsilon > 0$ , let  $\phi \in \Lambda_{\epsilon}(\mu^{\star})$ . By Corollary 1, it holds:

$$\max_{(\boldsymbol{\omega},\boldsymbol{\rho})\in\Omega\times\mathcal{P}} U^{\boldsymbol{\omega}\to\boldsymbol{\rho}}(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) \leq \sum_{\sigma=(I,a)\in\Sigma_{r}} \left[ \max_{\boldsymbol{\rho}_{\sigma}\in\mathcal{X}_{r,I}} U_{\sigma\to\boldsymbol{\rho}_{\sigma}}(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) - U(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}) \right]^{+}$$
$$\leq \sum_{\sigma=(I,a)\in\Sigma_{r}} \frac{\epsilon}{|\Sigma_{r}|} = \epsilon,$$

which implies that  $\phi \in \Phi_{\epsilon}^{\diamond}(\mu^{\star})$ . This concludes the proof.

**Theorem 2.** The BPSDM problem can be solved in polynomial time when the prior  $\mu^*$  is known.

662 *Proof.* It easy to check that the problem can be written as the following linear program:

$$\max_{\boldsymbol{\phi}\in\Lambda(\boldsymbol{\mu}^{\star})}F(\boldsymbol{\phi},\boldsymbol{\mu}^{\star}),$$

where the objective function is linear and  $\Lambda(\mu^*)$  is a polytope that can be represented by a polynomial number of linear inequalities, by Lemma 1.

### 665 E Proofs omitted from Section 5

**Theorem 3** (Impossibility of persuasiveness). There exists a constant  $\gamma \in (0, 1)$  such that no algorithm can guarantee to output a sequence  $\phi_1, \ldots, \phi_T$  of signaling schemes such that, with probability al least  $\gamma$ , all the signaling schemes  $\phi_t$  are persuasive.



Figure 4: Tree structure used in the proof of Theorem 3. Black round nodes are decision nodes  $\mathcal{H}_d$ . White round nodes are the chance nodes  $\mathcal{H}_c$ , while grey square nodes are the terminal nodes  $\mathcal{Z}$ .

*Proof.* We define two instances i and j of BPSDM problem based on the tree structure presented in Figure 4. In instance i, respectively j, the prior is defined as follows:

$$\mathfrak{i} := \begin{cases} \boldsymbol{\mu}^{\star}[(h_1, c)] = \frac{1}{2} + \epsilon \\ \boldsymbol{\mu}^{\star}[(h_1, d)] = \frac{1}{2} - \epsilon \\ \boldsymbol{\mu}^{\star}[(h_2, e)] = \frac{1}{2} - \epsilon \\ \boldsymbol{\mu}^{\star}[(h_2, f)] = \frac{1}{2} + \epsilon \end{cases}$$

671

$$\mathfrak{j} := \begin{cases} \boldsymbol{\mu}^{\star}[(h_1, c)] = \frac{1}{2} - \epsilon \\ \boldsymbol{\mu}^{\star}[(h_1, d)] = \frac{1}{2} + \epsilon \\ \boldsymbol{\mu}^{\star}[(h_2, e)] = \frac{1}{2} + \epsilon \\ \boldsymbol{\mu}^{\star}[(h_2, f)] = \frac{1}{2} - \epsilon \end{cases}$$

Moreover, for both instances  $u(z_1) = u(z_3) = 1$  and  $u(z_2) = u(z_4) = 0$ . A direct computation shows that, in instance i, it holds  $V_T^i = 2\epsilon \sum_{t=1}^T \phi_t[(h_0, b)]$ , while one can similarly compute that  $V_T^j = 2\epsilon \sum_{t=1}^T \phi_t[(h_0, a)]$ . Let  $\mathbb{P}^i$  and  $\mathbb{P}^j$  be the probability measures of instance i and j, respectively. Assume that  $\mathbb{P}^j[V_T^i \le 0] \ge 1 - \delta$ . Then, we know from the Pinsker inequality that:

$$\mathbb{P}^{\mathbf{i}}\left[\sum_{t=1}^{T} \boldsymbol{\phi}_{t}[(h_{0}, a)] \leq 0\right] \geq 1 - \sqrt{\frac{1}{2}\mathcal{K}(\mathbf{i}, \mathbf{j})} - \delta,$$

where  $\mathcal{K}(i, j)$  is the Kullback-Leibler divergence between instance i and j. By using the Kullback-Leibler decomposition (see, *e.g.*, [21] for more details), we can state that:

$$\mathcal{K}(\mathfrak{i},\mathfrak{j}) = 2T\mathcal{K}(B_{1/2+\epsilon}, B_{1/2-\epsilon}),$$

where  $\mathcal{K}(B_{1/2+\epsilon}, B_{1/2-\epsilon}) \leq 16\epsilon^2$  is the Kullback-Leibler divergence between a Bernoulli of parameter  $1/2 + \epsilon$  and one of parameter  $1/2 - \epsilon$ . Thus:

$$\mathbb{P}^{\mathbf{i}}\left[\sum_{t=1}^{T} \boldsymbol{\phi}_{t}[(h_{0}, a)] \leq 0\right] \geq 1 - 4\epsilon \sqrt{T} - \delta.$$

Moreover, in instance i, we have that  $V_T^i = 2\epsilon \sum_{t=1}^T \phi_t[(h_0, b)]$ , which implies:

$$\mathbb{P}^{\mathsf{i}}\left[V_T^{\mathsf{i}} \ge 2\epsilon T\right] \ge 1 - 4\epsilon \sqrt{T} - \delta.$$

By setting  $\epsilon = \frac{1}{16\sqrt{T}}$ , we have that:

$$\mathbb{P}^{\mathbf{i}}\left[V_T^{\mathbf{i}} \ge \frac{1}{8}\sqrt{T}\right] \ge 0.75 - \delta.$$

Thus, any algorithm that guarantees with high probability  $R_T^r \leq 0$  in instance j fails with high probability in instance i. This proves the claim.

#### F **Proofs omitted from Section 6** 684

Before presenting the proofs of the results in Section 6, we introduce some preliminary lemmas. 685

**Lemma 7.** Given any  $\phi \in \Phi$  and  $\mu, \mu' \in \mathcal{X}_c$ , if it is the case that  $\phi \in \Lambda_{\epsilon}(\mu)$  and  $\|\mu - \mu'\|_{\infty} \leq \gamma$ , then it holds that  $\phi \in \Lambda_{\epsilon'}(\mu')$  with  $\epsilon' = 2|\mathcal{Z}||\Sigma_r|\gamma + \epsilon$ . 686 687

*Proof.* For every  $(\sigma, \rho_{\sigma})$ -SPDP with  $\sigma = (I, a)$ , the following inequalities hold: 688

$$\begin{aligned} U_{\sigma \to \rho_{\sigma}}(\phi, \mu') &- U(\phi, \mu') \\ &= \sum_{\mathcal{Z}(I)} \phi[(h_{I}(z), a)] \rho_{\sigma}[\sigma_{r}(z)] \mu'[\sigma_{c}(z)] u(z) - \sum_{z \in \mathcal{Z}(\sigma)} \phi[\sigma_{s}(z)] \mu'[\sigma_{c}(z)] u(z) \\ &\leq \sum_{\mathcal{Z}(I)} \phi[(h_{I}(z), a)] \rho_{\sigma}[\sigma_{r}(z)] (\mu'[\sigma_{c}(z)] - \mu[\sigma_{c}(z)]) u(z) \\ &- \sum_{z \in \mathcal{Z}(\sigma)} \phi[\sigma_{s}(z)] (\mu'[\sigma_{c}(z)] - \mu[\sigma_{c}(z)]) u(z) + \frac{\epsilon}{|\Sigma_{r}|} \\ &\leq 2|\mathcal{Z}| ||\mu - \mu'||_{\infty} + \frac{\epsilon}{|\Sigma_{r}|} \leq 2|\mathcal{Z}|\gamma + \frac{\epsilon}{|\Sigma_{r}|}, \end{aligned}$$

where in the first inequality we added and subtracted the difference  $U_{\sigma \to \rho_{\sigma}}(\phi, \mu) - U(\phi, \mu)$  and 689 used the fact that  $\phi \in \Lambda_{\epsilon}(\mu)$ , while the second-to-last inequality follows from Hölder's inequality. 690 Since  $U_{\sigma \to \rho_{\sigma}}(\phi, \mu') - U(\phi, \mu') \leq 2|\mathcal{Z}|\gamma + \frac{\epsilon}{|\Sigma_r|} \coloneqq \frac{\epsilon'}{|\Sigma_r|}$  holds for every  $(\sigma, \rho_{\sigma})$ -SPDP, we have that  $\phi \in \Lambda_{\epsilon'}(\mu')$  with  $\epsilon' = |\mathcal{Z}||\Sigma_r|\gamma + \epsilon$ , concluding the proof. 691 692

**Lemma 8.** Given any  $\delta \in (0, 1)$ , Algorithm 1 guarantees that  $\mathbb{P}[\mathcal{E}] \ge 1 - \delta$ , where: 693

$$\mathcal{E} \coloneqq \{ \| \widehat{\boldsymbol{\mu}}_t - \boldsymbol{\mu}^{\star} \|_{\infty} \le \epsilon_t \quad \forall t \in [T] \},\$$

- and  $\epsilon_t$  is chosen according to Algorithm 1. 694
- *Proof.* Let  $\mathcal{B}_t(\delta)$  be defined as follows: 695

$$\mathcal{B}_t(\delta) \coloneqq \left\{ \boldsymbol{\mu} \left| \left| \boldsymbol{\mu}[\sigma] - \widehat{\boldsymbol{\mu}}_t[\sigma] \right| \le \sqrt{\frac{\log(2T|\Sigma_c|/\delta)}{2t}} \,\,\forall \sigma \in \Sigma_c \right\}.$$

Clearly,  $\mathbb{P}[\mathcal{E}] = \mathbb{P}[\mu^* \in \mathcal{B}_t(\delta) \ \forall t \in [T]]$ . By Hoeffding's inequality, we have that: 696

$$\mathbb{P}\left(|\boldsymbol{\mu}^{\star}[\sigma] - \widehat{\boldsymbol{\mu}}_{t}[\sigma]| \leq \sqrt{\frac{\log(2T|\Sigma_{c}|/\delta)}{2t}}\right) \geq 1 - \frac{\delta}{T|\Sigma_{c}|}$$

By a union bound over  $\sigma \in \Sigma_c$  and  $t \in [T]$ , we get that: 697

$$\mathbb{P}\left(|\boldsymbol{\mu}^{\star}[\sigma] - \widehat{\boldsymbol{\mu}}_{t}[\sigma]| \leq \sqrt{\frac{\log(2T|\Sigma_{c}|/\delta)}{2t}} \ \forall \sigma \in \Sigma_{c} \ \forall t \in [T]\right) \geq 1 - \delta.$$

- This concludes the proof of the lemma. 698
- **Lemma 9.** If the event  $\mathcal{E}$  holds, Algorithm 1 guarantees that  $\phi^* \in \Lambda_{\beta_t}(\widehat{\mu}_t)$  for all  $t \in [T]$ . 699
- *Proof.* By definition, we have that  $\phi^* \in \Lambda(\mu^*)$ . Moreover, since we conditioned on  $\mathcal{E}$ , we have that: 700  $\|\boldsymbol{\mu}^{\star} - \widehat{\boldsymbol{\mu}}_t\|_{\infty} \leq \epsilon_t \ \forall t \in [T].$
- Thus, we can exploit Lemma 7, which, by letting  $\beta_t \coloneqq 2|\mathcal{Z}||\Sigma_r|\epsilon_t$ , gives that  $\phi^* \in \Lambda_{\beta_t}(\widehat{\mu}_t)$ . 701
- **Lemma 10.** If the event  $\mathcal{E}$  holds, Algorithm 1 guarantees that  $\phi_t \in \Lambda_{2\beta_t}(\mu^*)$  for all  $t \in [T]$ . 702
- *Proof.* Given how Algorithm 1 works, we have that  $\phi_t \in \Lambda_{\beta_t}(\widehat{\mu}_t)$ . On the other hand, since we conditioned on the event  $\mathcal{E}$ , it must be the case that  $\|\boldsymbol{\mu}^* \widehat{\boldsymbol{\mu}}_t\| \leq \epsilon_t$  for all  $t \in [T]$ . Thus, by Lemma 7 703 704

**Theorem 4.** Given any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ , Algorithm 1 guarantees:

$$R_T = \mathcal{O}\left(|\mathcal{Z}|\sqrt{T\log(T|\Sigma_c|/\delta)}\right), \quad V_T = \mathcal{O}\left(|\Sigma_r||\mathcal{Z}|\sqrt{T\log(T|\Sigma_c|/\delta)}\right)$$

707 *Moreover, the algorithm runs in polynomial time.* 

*Proof.* First, we bound the computational complexity of the algorithm, then we separately analyze the sender's regret  $R_T$  and the receiver's regret  $V_T$ .

Complexity. With an argument analogous to the one used for the proof of Theorem 2, we have that the optimization problem solved by SELECTSTRATEGY() in Algorithm 1 is a polynomially-sized linear problem (Lemma 1). Hence, it can be solved in polynomial time.

713 Sender's regret. If the event  $\mathcal{E}$  holds, which happens with probability at least  $1 - \delta$ , then:

$$\boldsymbol{\mu}^{\star}[\sigma] - \boldsymbol{\epsilon}_t \leq \widehat{\boldsymbol{\mu}}_t[\sigma] \leq \boldsymbol{\mu}^{\star}[\sigma] + \boldsymbol{\epsilon}_t,$$

for every sequence  $\sigma \in \Sigma_c$  and round  $t \in [T]$ . This implies that, for every  $\phi \in \Phi$ , we have:

$$F(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) - |\mathcal{Z}|\epsilon_t \leq F(\boldsymbol{\phi}, \widehat{\boldsymbol{\mu}}_t) \leq F(\boldsymbol{\phi}, \boldsymbol{\mu}^{\star}) + |\mathcal{Z}|\epsilon_t.$$

- <sup>715</sup> Moreover, under the event  $\mathcal{E}$ , we have that  $\phi^* \in \Lambda_{\beta_t}(\widehat{\mu}_t)$  and, thus,  $F(\phi^*, \widehat{\mu}_t) \leq F(\phi_t, \widehat{\mu}_t)$  as  $\phi_t$  is
- computed by optimizing  $F(\cdot, \hat{\mu}_t)$  over  $\Lambda_{\beta_t}(\hat{\mu}_t)$ . By putting all the above results together, we get that, under event  $\mathcal{E}$ , the following holds:

$$F(\boldsymbol{\phi}^{\star},\boldsymbol{\mu}^{\star}) \leq F(\boldsymbol{\phi}^{\star},\widehat{\boldsymbol{\mu}}_{t}) + |\mathcal{Z}|\epsilon_{t} \leq F(\boldsymbol{\phi}_{t},\widehat{\boldsymbol{\mu}}_{t}) + |\mathcal{Z}|\epsilon_{t} \leq F(\boldsymbol{\phi}_{t},\boldsymbol{\mu}^{\star}) + 2|\mathcal{Z}|\epsilon_{t}.$$

By rearranging the terms, taking the sum over  $t \in [T]$ , and using  $\sum_{t=1}^{T} \frac{1}{\sqrt{t}} \le 2\sqrt{T}$ , we get:

$$R_T \coloneqq \sum_{t=1}^T \left( F(\boldsymbol{\phi}^*, \boldsymbol{\mu}^*) - F(\boldsymbol{\phi}_t, \boldsymbol{\mu}^*) \right) \le 2|\mathcal{Z}| \sum_{t=1}^T \epsilon_t \le 2|\mathcal{Z}| \sqrt{2\log(2T|\Sigma_c|/\delta)T},$$

which holds under the event  $\mathcal{E}$ , and, thus, with probability at least  $1 - \delta$ .

**Receiver's regret.** If the event  $\mathcal{E}$  holds, thanks to Lemma 10 we have that  $\phi_t \in \Lambda_{2\beta_t}(\mu^*)$ . Thus, by using Lemma 2, we can conclude that  $\phi_t \in \Phi_{2\beta_t}^{\diamond}(\mu^*)$ . This implies that, with probability at least  $1 - \delta$ , the following holds:

$$V_T \le 2\sum_{t=1}^T \beta_t \le 4|\Sigma_r||\mathcal{Z}|\sqrt{2\log(2T|\Sigma_c|/\delta)T},$$

view which concludes the proof.

### 724 G Proofs omitted from Section 7

Lemma 3. For every deterministic signaling scheme  $\pi \in \Pi$ , let

$$\Sigma_{\downarrow}(\boldsymbol{\pi}) \coloneqq \left\{ \sigma = (h, a) \in \Sigma_c \mid a \in A(h) \land \boldsymbol{\pi}[\sigma_s(h)] = 1 \right\}.$$

Then, during each round  $t \leq N$  of Algorithm 2, it holds  $\mathbb{E}[\mathbf{p}_t[\sigma]] = \boldsymbol{\mu}^{\star}[\sigma]$  for every  $\sigma \in \Sigma_{\downarrow}(\boldsymbol{\pi}_t)$ .

*Proof.* For any signaling scheme  $\phi \in \Phi$ , we have that the probability of reaching any node  $h \in \mathcal{H}_c$ 

during a round t < N (or, equivalently, that  $p_t[\sigma] = 1$  for some chance sequence  $\sigma = (h, a)$ ) is a Bernoulli with parameter  $\mu^*[\sigma]\phi[\sigma_s(h)]$ . Thus:

$$\mathbb{E}[\boldsymbol{p}_t[\sigma]] = \boldsymbol{\phi}_t[\sigma_s(h)]\boldsymbol{\mu}^{\star}[\sigma].$$

<sup>730</sup> If we consider any deterministic signaling scheme  $\pi \in \Pi$  and a chance sequence  $\sigma = (h, a) \in \Sigma_{\downarrow}(\pi)$ ,

we have that  $\phi_t[\sigma_s(h)] = 1$ , and, thus, the above equation simplifies to:

$$\mathbb{E}[\boldsymbol{p}_t[\sigma]] = \boldsymbol{\mu}^{\star}[\sigma]$$

via which concludes the proof.

**Lemma 11.** Given any  $\delta \in (0, 1)$ , Algorithm 2 guarantees that with probability at least  $1 - \delta/2$ : 733

$$\sum_{t=N+1}^{I} \sum_{z \in \mathcal{Z}} \epsilon_t[\sigma_c(z)] \phi_t[\sigma_s(z)] \le \sqrt{\log(4T|\Sigma_c|/\delta)|\Sigma_c|T} + |\mathcal{Z}|\sqrt{\log(2/\delta)T}$$

where the terms  $\epsilon_t[\sigma]$  for  $\sigma \in \Sigma_c$  and  $t \in [T]$  are defined according to Algorithm 2. 734

*Proof.* First, let us consider the deterministic signaling scheme  $\pi_t \in \Pi$  sampled by the algorithm 735

according to  $\phi_t$  at round  $t \in [T]$ . For convenience, in the following we report the definition of  $\epsilon_t[\sigma]$ 736 (according to Algorithm 2) for each  $\sigma \in \Sigma_c$  and  $t \in [T]$ : 737

$$\epsilon_t[\sigma] \coloneqq \sqrt{\frac{\log(4T|\Sigma_c|/\delta)}{2C_t[\sigma]}},$$

where  $C_t[\sigma]$  represents the number of rounds  $t' \leq t$  in which it is the case that  $\sigma \in \Sigma_{\downarrow}(\pi_{t'})$ . Then, 738 the following chain of inequalities holds: 739

$$\sum_{t=N+1}^{I} \sum_{z \in \mathcal{Z}} \epsilon_t[\sigma_c(z)] \boldsymbol{\pi}_t[\sigma_s(z)]$$
(9a)

$$=\sum_{t=N+1}^{T}\sum_{\substack{\sigma\in\Sigma_c:\\\exists z\in\mathcal{Z}:\sigma=\sigma_c(z)}}\left(\epsilon_t[\sigma]\sum_{\substack{\sigma'\in\Sigma_s:\\\exists z\in\mathcal{Z}:\sigma=\sigma_c(z)\wedge\sigma'=\sigma_s(z)}}\pi_t[\sigma']\right)$$
(9b)

$$\leq \sum_{t=N+1}^{1} \sum_{\sigma=(h,a)\in\Sigma_c} \epsilon_t[\sigma] \pi_t[\sigma_s(h)]$$
(9c)

$$= \sum_{\sigma=(h,a)\in\Sigma_c} \sum_{\substack{t\in[T]:\\t\geq N+1\wedge\pi_t[\sigma_s(h)]=1}} \epsilon_t[\sigma]$$
(9d)

$$=\sum_{\sigma\in\Sigma_c}\sum_{t=C_{N+1}[\sigma]}^{C_T[\sigma]}\sqrt{\frac{\log(4T|\Sigma_c|/\delta)}{2t}}$$
(9e)

$$\leq \sum_{\sigma \in \Sigma_c} \sqrt{\log(4T|\Sigma_c|/\delta)C_T[\sigma]} \tag{9f}$$

$$\leq \sqrt{\log(4T|\Sigma_c|/\delta)|\Sigma_c|T},\tag{9g}$$

where Equation (9c) follows by the definition of sequence-form signaling scheme of the sender, 740

Equation (9d) follows by exchanging the sums over  $\sigma \in \Sigma_c$  and  $t \in [T]$  and recalling that  $\pi_t$  is a 741 742

deterministic signaling scheme, Equation (9e) holds by definition of  $\epsilon$ , while Equation (9f) comes from  $\sum_{t=1}^{T} \frac{1}{\sqrt{t}} \leq 2\sqrt{T}$ . Finally, Equation (9g) follows from the Cauchy-Schwarz inequality. 743

Next, we provide a similar bound on  $\sum_{t=N+1}^{T} \sum_{z \in \mathcal{Z}} \epsilon_t[\sigma_c(z)] \phi_t[\sigma_s(z)]$ . We do this by exploiting the Azuma-Hoeffding inequality [14]. Indeed, we have that  $\mathbb{E}[\pi_t[\sigma]|\mathcal{F}_{t-1}] = \phi_t[\sigma]$ , where  $\mathcal{F}_{t-1}$  is 744 745 the filtration generated up to time t - 1 from the interaction between the algorithm and the BPSDM 746 problem. Thus, with probability at least  $1 - \delta/2$  the following holds: 747

$$\sum_{t=N+1}^{T} \sum_{z \in \mathcal{Z}} \epsilon_t[\sigma_c(z)] \phi_t[\sigma_s(z)] \le \sum_{t=N+1}^{T} \sum_{z \in \mathcal{Z}} \epsilon_t[\sigma_c(z)] \pi_t[\sigma_c(z)] + |\mathcal{Z}| \sqrt{\log(2/\delta)T}.$$

By combining the equation above with Equation (9f), we obtain: 748

$$\sum_{t=N+1}^{T} \sum_{z \in \mathcal{Z}} \epsilon_t [\sigma_c(z)] \phi_t[\sigma_s(z)] \le \sqrt{\log(4T|\Sigma_c|/\delta)|\Sigma_c|T} + |\mathcal{Z}|\sqrt{\log(2/\delta)T}$$

This concludes the proof. 749

- **Lemma 4.** Under the event  $\tilde{\mathcal{E}}$ , Algorithm 2 guarantees that  $\phi_t \in \Lambda_{2\beta_N}(\mu^*)$  at each round t > N.
- *Proof.* The proof is similar to the one of Lemma 10. If the event  $\tilde{\mathcal{E}}$  holds, then we have that:

$$\|\boldsymbol{\mu}^{\star} - \widehat{\boldsymbol{\mu}}_{N}\|_{\infty} \leq \max_{\sigma \in \Sigma} \epsilon_{t}[\sigma] \coloneqq \epsilon_{N}.$$

Moreover,  $\phi_t \in \Lambda_{\beta_N}(\widehat{\mu}_N)$  and we can use Lemma 7 to conclude that  $\phi_t \in \Lambda_{\beta_N+2\epsilon_N|\Sigma_r||\mathcal{Z}|}(\mu^*)$  for

all t > N. The proof follows from  $\beta_N \ge 2\epsilon_N |\mathcal{Z}||\Sigma_r|$ , since  $\epsilon_N \le \sqrt{\frac{\log(4T|\Sigma_c|/\delta)|\Sigma_c|}{2N}}$ .

Lemma 5. If the event  $\tilde{\mathcal{E}}$  holds, then, for every round t > N, it holds that  $\phi^* \in \Lambda_{\beta_N}(\hat{\mu}_t)$  and max $_{\mu \in \mathcal{C}_t(\delta)} F(\phi^*, \mu) \ge F(\phi^*, \mu^*)$ .

*Proof.* Since  $\phi^* \in \Lambda(\mu^*)$  and, under the event  $\tilde{\mathcal{E}}$ , it holds that:

$$\|\boldsymbol{\mu}^{\star} - \widehat{\boldsymbol{\mu}}_N\|_{\infty} \leq \max_{\sigma \in \Sigma} \epsilon_t[\sigma] \coloneqq \epsilon_N,$$

- ve can use Lemma 7 to conclude that  $\phi^* \in \Lambda_{2|\Sigma_c||\mathcal{Z}|\epsilon_N}(\widehat{\mu}_N)$ . The proof of the first statement
- is concluded by observing that  $\beta_N \geq 2|\Sigma_r||\mathcal{Z}|\epsilon_N$ , since  $\epsilon_N \leq \sqrt{\frac{\log(4T|\Sigma_c|/\delta)|\Sigma_c|}{2N}}$ . The second

statement directly follows from the observation that, under the event  $\tilde{\mathcal{E}}$ , it holds  $\mu^* \in \mathcal{C}_t(\delta)$ .

**Theorem 5.** Given any  $\delta \in (0, 1)$  and  $N \in [T]$ , Algorithm 2 guarantees:

$$R_T = \mathcal{O}\left(N + \sqrt{\log(T|\Sigma_c|/\delta)|\Sigma_c|T}\right) \quad and \quad V_T = \mathcal{O}\left(N + T|\mathcal{Z}|\sqrt{|\Sigma_c|\log(T|\Sigma_c|/\delta)/N}\right),$$

with probability at least  $1 - \delta$ . Moreover, the algorithm runs in polynomial time.

*Proof.* First, we bound the computational complexity of the algorithm, then we separately analyze the sender's regret  $R_T$  and the receiver's regret  $V_T$ .

**Complexity.** First, observe that  $F(\phi, \mu)$  is a linear function in  $\mu$  and it only has positive terms. Thus, for every  $\phi \in \Phi$ , the maximum over  $C_t(\delta)$  in the optimization problem solved during the second phase of the SELECTSTRATEGY() procedure is reached on the boundary of  $C_t(\delta)$ , so that larger entries of  $\mu$  provide larger objective values. Formally, we define:

$$\boldsymbol{\mu}_t \in \arg \max_{\boldsymbol{\mu} \in \mathcal{C}_t(\delta)} F(\boldsymbol{\phi}, \boldsymbol{\mu}),$$

which is independent of  $\phi$ . Then, for every  $\sigma \in \Sigma_c$ , we have that  $\mu_t[\sigma] = \hat{\mu}_t[\sigma] + \epsilon_t[\sigma]$ . Thus, we can compute the signaling scheme  $\phi_t$  with a linear program as follows:

$$\boldsymbol{\phi}_t \leftarrow \max_{\boldsymbol{\phi} \in \Lambda_{\beta_t}(\widehat{\boldsymbol{\mu}}_t)} F(\boldsymbol{\phi}, \boldsymbol{\mu}_t), \tag{10}$$

and, similarly to the proof of Theorem 4, we have that the optimization problem in Equation (10) is a polynomially-sized linear program by Lemma 1. Hence, it can be solved in polynomial time.

772 **Sender's regret.** Under the event  $\tilde{\mathcal{E}}$ , which happens with probability at least  $1 - \delta/2$ , we have that 773  $|\boldsymbol{\mu}^*[\sigma] - \hat{\boldsymbol{\mu}}_t[\sigma]| \leq \epsilon_t[\sigma]$  for all t > N. Thus,

$$\|\boldsymbol{\mu}^{\star} - \widehat{\boldsymbol{\mu}}_t\|_{\infty} \le \max_{\sigma \in \Sigma_c} \epsilon_t[\sigma] \coloneqq \epsilon_N.$$
(11)

Then, we can conclude that, under event  $\tilde{\mathcal{E}}$ , it holds  $\mu^*[\sigma] + 2\epsilon_t[\sigma] \ge \mu_t[\sigma]$ . This in turn implies:

$$F(\boldsymbol{\phi}_t, \boldsymbol{\mu}_t) \leq F(\boldsymbol{\phi}_t, \boldsymbol{\mu}^\star) + 2\sum_{z \in \mathcal{Z}} \epsilon_t[\sigma_c(z)] \boldsymbol{\phi}_t[\sigma_s(z)].$$

By Lemma 5, we have that, under event  $\tilde{\mathcal{E}}$ , it holds  $\phi^* \in \Lambda_{\beta_N}(\hat{\mu}_N)$ . Hence,  $F(\phi^*, \mu_t) \leq F(\phi_t, \mu_t)$ as  $\phi_t$  is computed as the optimum over  $\Lambda_{\beta_N}(\mu_t)$ . Moreover, by Lemma 5 we also have that  $F(\phi^*, \mu^*) \leq F(\phi^*, \mu_t)$ , which implies:

$$F(\boldsymbol{\phi}^{\star}, \boldsymbol{\mu}^{\star}) \leq F(\boldsymbol{\phi}^{\star}, \boldsymbol{\mu}_{t}) \leq F(\boldsymbol{\phi}_{t}, \boldsymbol{\mu}_{t}) \leq F(\boldsymbol{\phi}_{t}, \boldsymbol{\mu}^{\star}) + 2\sum_{z \in \mathcal{Z}} \epsilon_{t}[\sigma_{c}(z)]\boldsymbol{\phi}_{t}[\sigma_{s}(z)]$$

Then, we can decompose the sender's regret as:

$$R_T = \sum_{t=1}^N \left( F(\boldsymbol{\phi}^\star, \boldsymbol{\mu}^\star) - F(\boldsymbol{\phi}_t, \boldsymbol{\mu}^\star) \right) + \sum_{t=N+1}^T \left( F(\boldsymbol{\phi}^\star, \boldsymbol{\mu}^\star) - F(\boldsymbol{\phi}_t, \boldsymbol{\mu}^\star) \right)$$
$$\leq N + 2 \sum_{t=N+1}^T \sum_{z \in \mathcal{Z}} \epsilon_t [\sigma_c(z)] \boldsymbol{\phi}_t [\sigma_s(z)].$$

779 By using Lemma 11 and a union bound, we can conclude that with probability at least  $1 - \delta$ :

$$R_T \le N + 2\left(\sqrt{\log(4T|\Sigma_c|/\delta)|\Sigma_c|T} + |\mathcal{Z}|\sqrt{\log(2/\delta)T}\right).$$

**Receiver's regret.** By Lemma 4, under the event  $\tilde{\mathcal{E}}$ , we have that  $\phi_t \in \Lambda_{2\beta_N}(\mu^*)$  for all  $t \ge N$ . Moreover, by Lemma 2, it holds that  $\Lambda_{2\beta_N}(\mu^*) \subseteq \Phi_{2\beta_N}^{\diamond}(\mu^*)$ . Hence, with probability at least  $1 - \delta$ :

$$V_T \le N + 2T\beta_N = N + 4T|\mathcal{Z}||\Sigma_r|\sqrt{\frac{|\Sigma_c|\log(4T|\Sigma_c|/\delta)}{2N}}.$$

782 This concludes the proof.

**Theorem 6.** For any  $\alpha \in [1/2, 1]$ , there exists a constant  $\gamma \in (0, 1)$  such that no algorithm guarantees both  $R_T = o(T^{\alpha})$  and  $V_T = o(T^{1-\alpha/2})$  with probability greater than  $\gamma$ .

*Proof.* We define two instances i and j of a BPSDM problem whose tree structures are as in Figure 4. In both instances, we have that  $f(z_1) = f(z_2) = 0$  and  $f(z_3) = f(z_4) = 1$  for the sender, while  $u(z_1) = u(z_3) = 1$  and  $u(z_2) = u(z_4) = 0$  for the receiver. Moreover, in both instances we have that for the chance node  $h_1$  it holds  $\mu^*[(h_1, c)] = \mu^*[(h_1, d)] = 1/2$ . Instead, the two instances differ in the probabilities of chance node  $h_2$ , which are defined as follows:

$$\mathfrak{i} \coloneqq \begin{cases} \boldsymbol{\mu}^{\star}[(h_2, e)] = \frac{1}{2} - \epsilon \\ \boldsymbol{\mu}^{\star}[(h_2, f)] = \frac{1}{2} + \epsilon \end{cases}$$

790

$$\mathfrak{j} \coloneqq \begin{cases} \boldsymbol{\mu}^{\star}[(h_2, e)] = \frac{1}{2} + \epsilon \\ \boldsymbol{\mu}^{\star}[(h_2, f)] = \frac{1}{2} - \epsilon \end{cases}$$

<sup>791</sup> Simple calculations show that, in instance j, we have that the regret of the sender is:

$$R_T^{\mathbf{j}} = \sum_{t=1}^T \boldsymbol{\phi}_t[(h_0, a)]$$

- <sup>792</sup> Hence, if we require that (in high probability with respect to the measure  $\mathbb{P}^{j}$  of instance j) the sender's
- regret is smaller than a threshold K, then:

$$\mathbb{P}^{\mathsf{j}}\left[\sum_{t=1}^{T} \boldsymbol{\phi}_t[(h_1, a)] \le K\right] \ge 1 - \delta.$$

794 The Pinsker's inequality states that:

$$\mathbb{P}^{\mathfrak{i}}\left[\sum_{t=1}^{T} \boldsymbol{\phi}_{t}[(h_{1}, a)] \leq K\right] \geq 1 - \delta - \sqrt{\frac{1}{2}\mathcal{K}(\mathfrak{j}, \mathfrak{i})},$$

where  $\mathcal{K}(j, i)$  is the Kullback-Leibler divergence between instance j and instance i. By the well-known decomposition theorem of the divergence, we know that:

$$\mathcal{K}(\mathbf{j},\mathbf{i}) = \mathbb{E}^{\mathbf{j}} \left[ \sum_{t=1}^{T} \boldsymbol{\phi}_t[(h_1, a)] \right] \mathcal{K}(B_{1/2+\epsilon}, B_{1/2-\epsilon}) \le 16\epsilon^2 \mathbb{E}^{\mathbf{j}} \left[ \sum_{t=1}^{T} \boldsymbol{\phi}_t[(h_1, a)] \right],$$

where  $\mathcal{K}(B_{1/2+\epsilon}, B_{1/2-\epsilon})$  is the Kullback-Leibler divergence between two Bernoulli random variable with parameter  $1/2 + \epsilon$  and  $1/2 - \epsilon$ . Now, we can upper bound  $\mathbb{E}^{j}\left[\sum_{t=1}^{T} \phi_{t}[(h_{1}, a)]\right]$  in terms of the probability  $\mathbb{P}^{j}$  with the reverse Markov inequality, as follows:

$$\mathbb{E}^{\mathsf{j}}\left[\sum_{t=1}^{T} \boldsymbol{\phi}_{t}[(h_{1}, a)]\right] \leq \mathbb{P}^{\mathsf{j}}\left[\sum_{t=1}^{T} \boldsymbol{\phi}_{t}[(h_{1}, a)] \geq K\right] (T - K) + K$$
$$\leq \delta(T - K) + K.$$

800 Thus, we can conclude that:

$$\mathbb{P}^{i}\left[\sum_{t=1}^{T} \boldsymbol{\phi}_{t}[(h_{1}, a)] \leq K\right] \geq 1 - \delta - 2\epsilon \sqrt{2(\delta(T - K) + K)}.$$
(12)

Now, we consider the receiver's regret in instance i, which can be computed as:

$$V_T^{\mathbf{i}} = \epsilon \sum_{t=1}^T \boldsymbol{\phi}_t[(h_0, b)].$$

802 This, together with Equation (12), allows us to conclude that:

$$\mathbb{P}^{\mathbf{i}}\left[V_T^{\mathbf{i}} \ge \epsilon(T-K)\right] \ge 1 - \delta - 2\epsilon\sqrt{2(\delta(T-K)+K)}$$

By setting  $K = \frac{T^{\alpha}}{8}$  and  $\epsilon = \frac{T^{-\alpha/2}}{8}$ , we can conclude that if

$$\mathbb{P}^{\mathsf{j}}\left[\sum_{t=1}^{T} \boldsymbol{\phi}_t[(h_0, a)] \le \frac{T^{\alpha}}{8}\right] \ge 1 - \delta,$$

804 then

$$\mathbb{P}^{i}\left[V_{T}^{i} \geq \frac{T^{1-\alpha/2}}{16}\right] \geq 1 - \frac{\sqrt{2}}{16} - \delta \geq 0.91 - \delta,$$

where we used that  $\frac{T^{1-\alpha/2}}{8} - \frac{T^{\alpha/2}}{64} \ge \frac{T^{1-\alpha/2}}{16}$  for  $T \ge 1$  and that we can assume  $\delta \le \frac{T^{\alpha-1}}{4}$ .  $\Box$