



Colonel Blotto and Hide-and-Seek Games as Path Planning Problems with Side Observations

Dong Quan Vu, Patrick Loiseau, Alonso Silva, Long Tran-Thanh

► To cite this version:

Dong Quan Vu, Patrick Loiseau, Alonso Silva, Long Tran-Thanh. Colonel Blotto and Hide-and-Seek Games as Path Planning Problems with Side Observations. 2019. hal-02139519v1

HAL Id: hal-02139519

<https://hal.science/hal-02139519v1>

Preprint submitted on 27 May 2019 (v1), last revised 28 May 2019 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Colonel Blotto Games and Hide-and-Seek Games as Path Planning Problems with Side Observations

Dong Quan Vu
Nokia Bell Labs France
AAAIRD department

Patrick Loiseau
Univ. Grenoble Alpes,
Inria, CNRS Grenoble INP LIG
& MPI-SWS

Alonso Silva
Safran Tech
Signal and Information
Technologies

Long Tran-Thanh
University of Southampton
Electronic and Computer
Science Department

Abstract

Resource allocation games such as the famous Colonel Blotto (CB) and Hide-and-Seek (HS) games are often used to model a large variety of practical problems, but only in their one-shot versions. Indeed, due to their extremely large strategy space, it remains an open question how one can efficiently learn in these games. In this work, we show that the online CB and HS games can be cast as path planning problems with side-observations (SOPPP): at each stage, a learner chooses a path on a directed acyclic graph and suffers the sum of losses that are adversarially assigned to the corresponding edges; and she then receives semi-bandit feedback with side-observations (i.e., she observes the losses on the chosen edges plus some others). Then, we propose a novel algorithm, EXP3-OE, the first-of-its-kind with guaranteed efficient running time for SOPPP without requiring any auxiliary oracle. We provide an expected-regret bound of EXP3-OE in SOPPP matching the order of the best benchmark in the literature. Moreover, we introduce additional assumptions on the observability model under which we can further improve the regret bounds of EXP3-OE. We illustrate the benefit of using EXP3-OE in SOPPP by applying it to the online CB and HS games.

1 Introduction

Resource allocation games have been studied profoundly in the literature and showed to be very useful to model many practical situations, including online decision problems, see e.g. [10, 12, 24, 40]. In particular, two of the most renowned are the Colonel Blotto game (henceforth, CB game) and the Hide-and-Seek game (henceforth, HS game). In the (one-shot) *CB game*, two players, each with a fixed amount of budget, simultaneously allocate their (indivisible) resources on $n \in \mathbb{N}$ battlefields, each player's payoff is the aggregate of the values of battlefields where she has a higher allocation. The scope of applications of the CB games includes a variety of problems; for instance, in security (e.g., [15, 32]) where resources correspond to security forces, in politics (e.g., [25, 30]) for allocating budget to attract voters, and in advertisement (e.g., [27, 28]) for distributing the broadcasting time. On the other hand, in the (one-shot) *HS game*, a seeker chooses n among k locations ($n < k$) to search for a hider, who randomly chooses to hide in one of the k locations. The seeker's payoff is the probability that she finds the hider and the hider's payoff is the probability that she successfully escape the seeker's pursuit. Several variants of the HS games are used to model surveillance situations [8, 9], anti-jamming problems in telecommunications [29, 37, 38], vehicles control [16, 34], etc.

Both the CB games and the HS games have a long-standing history (originated in 1921 [11] and 1953 [35], respectively); however, the results achieved so-far in these games are mostly limited to their one-shot and full-information version (see e.g., [6, 18, 30, 32, 36] for CB games and [19, 21, 39] for HS games). On the contrary, in most of the applications (e.g., web security, advertising, telecommunications), a more natural setting is to consider the case where the game is played repeatedly and players have access only to incomplete information at each stage. In this setting, players are often required to sequentially learn the game on-the-fly and adjust the trade-off between exploiting known information and exploring to gain new information. Thus, this work focuses on the following sequential learning problem: at each stage, a learner plays a CB game (resp. HS game); at the end of the stage, she receives limited feedback that is the gain she obtains from each battlefield (resp. the hider’s escape probability corresponding to the chosen locations); and her objective is to maximize her cumulative payoffs. A formal definition of these problems is given in Section 4; hereinafter, we reuse the term CB game and HS game to refer to this sequential learning version of the games. The main challenge in those games is that their strategy space is exponential in the natural parameters (e.g., number of troops and battlefields in the CB game, number of locations in the HS game); hence how to efficiently learn in these games is an open question.

Our **first contribution** towards solving this open question is to show that the CB and HS games can be cast as *Path Planning Problems* (henceforth, PPP), one of the most well-studied instances of the *Online Combinatorial Optimization* framework (henceforth, OCOMB; see [14] for a survey). In PPPs, given a graph with E edges, at each stage, a learner chooses a path; then a loss in $[0, 1]$ is adversarially chosen for each edge and the learner suffers the aggregate of edges’ losses belonging to her chosen path. The learner’s goal is to minimize regret.¹ The information that the learner receives in the CB and HS games as described above straightforwardly corresponds to the so-called *semi-bandit* feedback setting of PPPs, i.e., at the end of each stage, the learner observes the edges’ losses belonging to her chosen path. However, the specific structure of the considered games also allows the learner to deduce (without any extra cost) from the semi-bandit feedback the losses of some of the other edges that may not belong to the chosen path; these are called *side-observations*. Henceforth, we will use the term SOPPP to refer to this PPP under semi-bandit feedback with side-observations.

SOPPP is a special case of OCOMB with side-observations (henceforth, SOCOMB) studied by [23] and, following their approach, we will use *observation graphs*² (defined in Section 2) to capture the learner’s observability. In [23], the authors focus on the class of Follow-the-Perturbed-Leader (FPL) algorithms (originated from [22]) and propose an algorithm named FPL-IX for SOCOMB, which could be applied directly to SOPPP. However, this faces two main problems: (i) the efficiency of FPL-IX is only guaranteed with high-probability (as it depends on the geometric sampling technique) and (ii) it requires that there exists an efficient oracle that solves an optimization problem at each stage—both of which are incompatible with our goal of learning in the CB and HS games.

In this paper, we focus instead on another prominent class of OCOMB algorithms, called EXP3 [5, 17]. Then, our **second contribution** is to propose an algorithm for SOPPP that solves both of the aforementioned issues and provides good regret guarantees. In more details, this contribution is three-fold: (i) We propose a *novel algorithm*, EXP3-OE, that is applicable to any instance of SOPPP. Importantly, EXP3-OE is always guaranteed to run efficiently (i.e., in polynomial time in terms of the number of edges of the graph in SOPPP) without the need of any auxiliary oracle; (ii) We prove that EXP3-OE guarantees an upper-bound on the expected regret matching in order with the best benchmark in the literature (the FPL-IX algorithm). We also prove further improvements under additional assumptions on the observation graphs that have been so-far ignored in the literature; (iii) We demonstrate the benefit of using the EXP3-OE algorithm in the CB and HS games.

Our EXP3-OE algorithm is based on the EXP3-IX algorithm [23]. However, EXP3-IX has a very inefficient running time in SOCOMB (and particularly in SOPPP) and thus, it is only analyzed by [23] in the trivial cases of SOCOMB involving only actions with L1-norm that equals to 1 (corresponding to SOPPP with graphs where all paths have length 1)—the existence of an efficient implementation of EXP3-type algorithms in SOCOMB is left as an open question in [23]. We address this question in the particular case of SOPPP as follows. We introduce two main major updates in EXP3-OE. First, unlike EXP3-IX that uses adaptive implicit exploration scheme, we assume that the time horizon is known³

¹The regret is the difference between the learner’s cumulative loss and that of the best action in hindsight.

²The observation graphs, proposed in [23] and used here for SOPPP, extend the side-observations model for multi-armed bandits problems studied by [1, 2, 26]. Indeed, they capture side-observations between edges whereas the side-observations model considered in [1, 2, 26] is between actions (i.e., paths in PPPs).

³If T is unknown, we can use the doubling trick (see [4, 7]) to get similar results.

in advance and fix an implicit exploration parameter in the loss estimator of EXP3-OE. This change reduces the computations and leads to a different parameters tuning scheme with improved regret bounds compared to EXP3-IX. Second (and the main reason that makes EXP3-OE significantly more efficient than EXP3-IX), we use a novel loss estimator, which enables us to efficiently compute it based on a dynamic-programming technique, called *weight pushing*. Note that while *weight pushing* has been used for efficiently sampling paths from exponentially-updated weights in several variants of EXP3 (e.g., [20, 31, 33]), the way we apply it to compute the loss estimator is novel and non-trivial. Finally, note that the SOPPP model (and thus, our proposed EXP3-OE algorithm) can be applied into many problems beyond the considered games, e.g., auctions, recommendation systems.

Throughout the paper, we use bold symbols to denote vectors, e.g., $\mathbf{x} \in \mathbb{R}^n$, and $\mathbf{x}(i)$ to denote the i -th element. For any $m \geq 1$, the set $\{1, 2, \dots, m\}$ is denoted by $[m]$ and the indicator function of a set A is denoted by \mathbb{I}_A . For graphs, we write either $e \in \mathbf{p}$ or $\mathbf{p} \ni e$ to refer that an edge e belongs to a path \mathbf{p} . For the sake of conciseness, we present first our second contribution on the SOPPP in general and we then return in Section 4 to our first contribution relating to the CB and HS games.

2 Path Planning Problems with Side-Observations (SOPPP) Formulation

As discussed in Section 1, motivated by the CB and HS games, we focus on the path planning problem with semi-bandit and side-observations feedback (SOPPP) and design an EXP3-type algorithm that always runs efficiently in SOPPP. To do this, we first formally define the SOPPP model as follows.

SOPPP model. Consider a directed acyclic graph (henceforth, DAG), denoted by G , whose set of vertices and set of edges are respectively denoted by \mathcal{V} and \mathcal{E} . Let $V := |\mathcal{V}| \geq 2$ and $E := |\mathcal{E}| \geq 1$; there are two special vertices, a source and a destination, that are respectively called s and d . We denote by \mathcal{P} the set of all *paths* starting from s and ending at d . Each path $\mathbf{p} \in \mathcal{P}$ corresponds to a vector in $\{0, 1\}^E$ (thus, $\mathcal{P} \subset \{0, 1\}^E$) where $\mathbf{p}(e) = 1$ if and only if edge $e \in \mathcal{E}$ belongs to \mathbf{p} . Let n be the length of the longest path in \mathcal{P} , that is $\|\mathbf{p}\|_1 \leq n, \forall \mathbf{p} \in \mathcal{P}$. Given a time horizon $T \in \mathbb{N}$, at each (discrete) stage $t \in [T]$, a *learner* chooses a path $\tilde{\mathbf{p}}_t \in \mathcal{P}$. Then, a *loss vector* $\ell_t \in [0, 1]^E$ is secretly and adversarially chosen (oblivious from the learner's decisions). Each element $\ell_t(e)$ corresponds to the scalar loss embedded on the edge $e \in \mathcal{E}$. The learner's incurred loss is $L_t(\tilde{\mathbf{p}}_t) = (\tilde{\mathbf{p}}_t)^\top \ell_t = \sum_{e \in \tilde{\mathbf{p}}_t} \ell_t(e)$, i.e., the sum of the losses from all the edges belonging to $\tilde{\mathbf{p}}_t$. The learner's feedback at stage t after choosing $\tilde{\mathbf{p}}_t$ is presented as follows. First, she receives a *semi-bandit* feedback, that is, she observes all the edges' losses $\ell_t(e)$, for any e belonging to the chosen path $\tilde{\mathbf{p}}_t$. Additionally, each edge $e \in \tilde{\mathbf{p}}_t$ may reveal the losses on several other edges. To represent these *side-observations* at time t , we consider a graph, denoted G_t^O , containing E vertices. Each vertex v_e of G_t^O corresponds to an edge $e \in \mathcal{E}$ of the graph G . There exists a directed edge from a vertex v_e to a vertex $v_{e'}$ in G_t^O if, by observing the edge loss $\ell_t(e)$, the learner can also deduce the edge loss $\ell_t(e')$ (we also denote this by $e \rightarrow e'$ and say that the edge e reveals the edge e'). The objective of the learner is to minimize the cumulative *expected regret*, defined as $R_T := \mathbb{E} \left[\sum_{t \in [T]} L(\tilde{\mathbf{p}}_t) \right] - \min_{\mathbf{p}^* \in \mathcal{P}} \sum_{t \in [T]} L(\mathbf{p}^*)$.

Hereinafter, in places where there is no ambiguity, we use the term *path* to refer to a path in \mathcal{P} and the term *observation graphs* to refer to G_t^O . In general, these observation graphs can depend on the decisions of both the learner and the adversary. On the other hand, all vertices in G_t^O always have self-loops. In the case where none among $G_t^O, t \in [T]$ contains any other edge than these self-loops, no side-observation is allowed and the problem is reduced to the classical semi-bandit setting. If all $G_t^O, t \in [T]$ are complete graphs, SOPPP corresponds to the full-information PPPs. In this work, we focus on considering the *uninformed setting*, i.e., the learner observes G_t^O only after making a decision at time t . On the other hand, let us introduce two new notations:

$$\mathbb{O}_t(e) := \{\mathbf{p} \in \mathcal{P} : \exists e' \in \mathbf{p}, e' \rightarrow e\}, \forall e \in \mathcal{E} \text{ and } \mathbb{O}_t(\mathbf{p}) := \{e \in \mathcal{E} : \exists e' \in \mathbf{p}, e' \rightarrow e\}, \forall \mathbf{p} \in \mathcal{P}.$$

Intuitively, $\mathbb{O}_t(e)$ is the set of all paths that, if chosen, reveal the loss on the edge e and $\mathbb{O}_t(\mathbf{p})$ is the set of all edges whose losses are revealed if the path \mathbf{p} is chosen. Trivially, $\mathbf{p} \in \mathbb{O}_t(e) \Leftrightarrow e \in \mathbb{O}_t(\mathbf{p})$. Moreover, due to the semi-bandit feedback, if $\mathbf{p}^* \ni e^*$, then $\mathbf{p}^* \in \mathbb{O}_t(e^*)$ and $e^* \in \mathbb{O}_t(\mathbf{p}^*)$. Apart from the results for general observation graphs, in this work, we additionally present several results under two particular assumptions, satisfied by some instances in practice (e.g., the CB and HS games), that provide more refined regret bounds compared to cases that were considered in [23]: (i) *symmetric* observation graphs where for each edge from v_e to $v_{e'}$, there also exists an edge from $v_{e'}$ to v_e (i.e., if $e \rightarrow e'$ then $e' \rightarrow e$); i.e., G_t^O is an undirected graph; (ii) observation graphs under

the following *assumption* (A0) that requires that if two edges belong to a path in G , then they cannot simultaneously reveal the loss of another edge.

Assumption (A0): For any $e \in \mathcal{E}$, if $e' \rightarrow e$ and $e'' \rightarrow e$, then $\nexists \mathbf{p} \in \mathcal{P} : \mathbf{p} \ni e', \mathbf{p} \ni e''$.

3 EXP3-OE - An Efficient Algorithm for the SOPPP

In this section, we present a new algorithm for SOPPP, called EXP3-OE (OE stands for Observable Edges), whose pseudo-code is given by Algorithm 1. The guarantees on the expected regret of EXP3-OE in SOPPP is analyzed in Section 3.2. More importantly, EXP3-OE always runs efficiently in polynomial time in terms of the number of edges of G ; this is discussed in Section 3.1.

As an EXP3-type algorithm, EXP3-OE relies on the average weights sampling where at stage t we update the weight $w_t(e)$ on each edge e by the exponential rule (line 9). For each path \mathbf{p} , we denote the path weight $w_t(\mathbf{p}) := \prod_{e \in \mathbf{p}} w_t(e)$ and define the following normalized terms, according to which a path $\tilde{\mathbf{p}}_t$ is sampled at each stage t (see line 5) of the EXP3-OE algorithm:

$$d_t(\mathbf{p}) := \prod_{e \in \mathbf{p}} w_t(e) / \sum_{\mathbf{p}' \in \mathcal{P}} \prod_{e' \in \mathbf{p}'} w_t(e') = w_t(\mathbf{p}) / \sum_{\mathbf{p}' \in \mathcal{P}} w_t(\mathbf{p}'), \forall \mathbf{p} \in \mathcal{P}. \quad (1)$$

Compared to other instances of the EXP3-type algorithms, EXP3-OE has two major differences. First, at each stage t , the loss of each edge e is estimated by $\hat{\ell}_t(e)$ (line 8) based on the term $q_t(e)$ and a parameter β . Intuitively, $q_t(e)$ is the probability that the loss on the edge e is revealed from playing the chosen path at t . On the other hand, the implicit exploration parameter β added to the denominator allows us to “pretend to explore” in EXP3-OE without knowing the observation graph G_t^O before making the decision at stage t (the uninformed setting). Unlike the standard EXP3 algorithm, the loss estimator used in EXP3-OE is *biased*, that is

Algorithm 1 EXP3-OE Algorithm for SOPPP.

-
- 1: **Input:** $T, \eta, \beta > 0$, graph G .
 - 2: Initialize $w_1(e) := 1, \forall e \in \mathcal{E}$.
 - 3: **for** $t = 1$ **to** T **do**
 - 4: Loss vector ℓ_t is chosen adversarially (unobserved).
 - 5: Sample a path $\tilde{\mathbf{p}}_t$ according to $d_t(\tilde{\mathbf{p}}_t)$ by Algorithm 4 (Appendix A).
 - 6: Suffer the loss $L_t(\tilde{\mathbf{p}}_t) = \sum_{e \in \tilde{\mathbf{p}}_t} \ell_t(e)$.
 - 7: Observation graph G_t^O is generated and $\ell_t(e), \forall e \in \mathbb{O}_t(\tilde{\mathbf{p}}_t)$ are observed.
 - 8: $\hat{\ell}_t(e) := \frac{\ell_t(e)}{q_t(e) + \beta} \mathbb{I}_{\{e \in \mathbb{O}_t(\tilde{\mathbf{p}}_t)\}}, \forall e \in \mathcal{E}$, where $q_t(e) := \sum_{\mathbf{p} \in \mathbb{O}_t(e)} d_t(\mathbf{p})$ is computed by Algorithm 2.
 - 9: Update weights $w_{t+1}(e) := w_t(e) \cdot \exp(-\eta \hat{\ell}_t(e))$.
 - 10: **end for**
-

$$\mathbb{E}_t \left[\hat{\ell}_t(e) \right] = \sum_{\tilde{\mathbf{p}} \in \mathcal{P}} d_t(\tilde{\mathbf{p}}) \frac{\ell_t(e)}{q_t(e) + \beta} \mathbb{I}_{\{e \in \mathbb{O}_t(\tilde{\mathbf{p}})\}} = \sum_{\tilde{\mathbf{p}} \in \mathbb{O}_t(e)} d_t(\tilde{\mathbf{p}}) \frac{\ell_t(e)}{\sum_{\mathbf{p} \in \mathbb{O}_t(e)} d_t(\mathbf{p}) + \beta} \leq \ell_t(e), \forall e \in \mathcal{E}. \quad (2)$$

Here, \mathbb{E}_t denotes the expectation w.r.t. the randomness of choosing a path at stage t . Second, unlike standard EXP3 algorithms that keep track and update on the weight of each path, the weight pushing technique is applied at line 5 (via Algorithm 4 in Appendix A) and line 8 (via Algorithm 2 in Section 3.1) where we work with edges weights instead of paths weights (recall that $E \ll |\mathcal{P}|$).

3.1 Running Time Efficiency of the EXP3-OE Algorithm

We recall that in order to efficiently sample a path according to $d_t(\mathbf{p}), \mathbf{p} \in \mathcal{P}$, following the literature, it is useful to compute the terms $H_t(s, u) := \sum_{\mathbf{p} \in \mathcal{P}_{s,u}} \prod_{e \in \mathbf{p}} w_t(e)$ and $H_t(u, d) := \sum_{\mathbf{p} \in \mathcal{P}_{u,d}} \prod_{e \in \mathbf{p}} w_t(e)$ for any vertex u in G . Intuitively, $H_t(u, v)$ is the aggregate weight of all paths from vertex u to vertex v at stage t . Then, a path in G is sampled sequentially edge-by-edge based on these terms H_t . The collection of the computations described above is often referred to as weight pushing, that can be done in $\mathcal{O}(E)$ by exploiting the structure of the graph. We rewrite this step formally in Appendix A.

The final non-trivial step to efficiently implement EXP3-OE is to compute $q_t(e)$, the probability that an edge e is revealed at stage t , needed in line 8. We note that $q_t(e)$ is the sum of $|\mathbb{O}_t(e)| = \mathcal{O}(|\mathcal{P}|)$ terms; therefore, a direct computation is inefficient while a naive application of the weight pushing technique can easily lead to errors. To compute $q_t(e)$, we propose Algorithm 2, a non-straightforward application of weight pushing, in which we consecutively consider all the edges $e' \in \mathfrak{R}_t(e) := \{e' \in \mathcal{E} : e' \rightarrow e\}$. Then, we take the sum of the terms $d_t(\mathbf{p})$ of the paths \mathbf{p} going through e' by the weight pushing

technique while making sure that each of these terms $d_t(\mathbf{p})$ is only included one time, even if \mathbf{p} has more than one edge revealing e (this is a non-trivial step). In Algorithm 2, we denote by $C(u)$ the set of the direct successors of any vertex $u \in \mathcal{V}$. A proof that Algorithm 2 outputs exactly $q_t(e)$ as defined in line 8 of Algorithm 1 can be found in Appendix B. Algorithm 2 runs in $\mathcal{O}(|\mathfrak{R}_t(e)|E)$ time; therefore, line 8 of Algorithm 1 can be done in at most $\mathcal{O}(E^3)$ time. In conclusion, the EXP3-OE algorithm runs in at most $\mathcal{O}(E^3T)$ time, this guarantee works even for the worst-case scenario. For comparison, the running time of FPL-IX proposed by [23] is $\mathcal{O}(E|\mathcal{V}|^2T)$ in expectation if we choose Dijkstra’s algorithm to be the optimization oracle at each stage. On the other hand, with the chosen parameters in [23], we can deduce that FPL-IX achieves the running time in⁴ $\tilde{\mathcal{O}}(n^{1/2}E^{3/2}\ln(E/\delta)T^{3/2})$ with a probability at least $1 - \delta$ for an arbitrary $\delta > 0$. That is, FPL-IX is not guaranteed to have efficient running time in all cases.

Algorithm 2 Compute $q_t(e)$ of an edge e at stage t .

```

1: Input:  $e \in \mathbb{O}_t(\tilde{\mathbf{p}}_t)$ , set  $\mathfrak{R}_t(e)$  and  $w_t(\bar{e}), \forall \bar{e} \in \mathcal{E}$ .
2: Initialize  $\bar{w}(\bar{e}) := w_t(\bar{e}), \forall \bar{e} \in \mathcal{E}$  and  $q_t(e) := 0$ .
3: Compute  $H^*(s, d)$  by Algorithm 3 (Appendix A) with input  $\{w_t(\bar{e}), \bar{e} \in \mathcal{E}\}$ .
4: for  $e' \in \mathfrak{R}_t(e)$  do
5:   Compute  $H(s, u), H(u, d), \forall u \in \mathcal{V}$  by Algorithm 3 with input  $\{\bar{w}(\bar{e}), \bar{e} \in \mathcal{E}\}$ .
6:    $K(e') := H(s, u_{e'}) \cdot w(e') \cdot H(v_{e'}, d)$  where edge  $e'$  goes from  $u_{e'}$  to  $v_{e'} \in C(u_{e'})$ .
7:    $q_t(e) := q_t(e) + K(e')/H^*(s, d)$ .
8:   Update  $\bar{w}(e') = 0$ .
9: end for
10: Output:  $q_t(e)$ .
```

3.2 Performance of the EXP3-OE Algorithm

In this section, we present an upper-bound of the expected regret achieved by the EXP3-OE algorithm in the SOPPP. For the sake of brevity, with $d_t(\mathbf{p})$ defined in (1), for any $t \in [T]$ and $e \in \mathcal{E}$, we denote:

$$r_t(e) := \sum_{\mathbf{p} \ni e} d_t(\mathbf{p}) \text{ and } Q_t := \sum_{e \in \mathcal{E}} r_t(e)/(q_t(e) + \beta). \quad (3)$$

Intuitively, $r_t(e)$ is the probability that the chosen path at stage t contains an edge e and Q_t is the summation over all the edges of the ratio of this quantity and the probability that the loss of an edge is revealed (plus β). We can bound the expected regret with this key term Q_t .

Theorem 3.1. *The expected regret of the EXP3-OE algorithm in the SOPPP satisfies:*

$$R_T \leq \ln(|\mathcal{P}|)/\eta + [\beta + (n \cdot \eta)/2] \cdot \sum_{t \in [T]} Q_t. \quad (4)$$

The proof of Theorem 3.1 is given in Appendix C and has an approach similar to [2, 13] with several necessary adjustments to handle the new biased loss estimator in EXP3-OE. To see the relationship between the structure of the side-observations of the learner and the bound of the expected regret, we look for the upper-bounds of Q_t in terms of the observation graphs’ parameters. Let α_t be the independence number⁵ of G_t^O , we have the following statement.

Theorem 3.2. *Let us denote $M := \lceil 2E^2/\beta \rceil$, $N_t := \ln\left(1 + \frac{M+E}{\alpha_t}\right)$ and $K_t := \ln\left(1 + \frac{nM+E}{\alpha_t}\right)$, Upper-bounds of Q_t in different cases of G_t^O are given in the following table:*

	SATISFIES (A0)	NOT SATISFIES (A0)
SYMMETRIC	α_t	$n\alpha_t$
NON-SYMMETRIC	$1+2\alpha_t N_t$	$2n(1+\alpha_t K_t)$

A proof of this theorem is given in Appendix E. The main idea of this proof is based on several graph theoretical lemmas that are extracted from [2, 23, 26]. These lemmas establish the relationship between the independence number of a graph and the ratios of the weights on the graph’s vertices that have similar forms to the key-term Q_t . The case where observation graphs are non-symmetric

⁴The notation $\tilde{\mathcal{O}}$ is a version of the big-O asymptotic notation that ignores the logarithmic terms.

⁵The independence number of a directed graph is computed while ignoring the direction of the edges.

and do not satisfy assumption (A0) is the most general setting. Moreover, as showed in Theorem 3.2, the bounds of Q_t are improved if the observation graphs satisfy either the symmetry condition or assumption (A0). Intuitively, given the same independence numbers, a symmetric observation graph gives the learner more information than a non-symmetric one; thus, it may yield a better bound on Q_t and the expected regret. On the other hand, assumption (A0) is a technical assumption that allows the use of different techniques in the proofs to obtain better bounds. These cases have not been analyzed in the literature while they are satisfied by several practical situations, including the CB and HS games (see Section 4).

Finally, we give results on the order of the upper-bounds of the expected regret, obtained by the EXP3-OE algorithm, presented as a corollary of Theorems 3.1 and 3.2.

Corollary 3.3. *In SOPPP, let α be an upper bound of $\alpha_t, \forall t \in [T]$. With appropriate choices of the parameters η and β , the expected regret of the EXP3-OE algorithm is:*

- (i) $R_T \leq \tilde{O}(n\sqrt{T\alpha\ln(|\mathcal{P}|)})$ in the general cases.
- (ii) $R_T \leq \tilde{O}(\sqrt{nT\alpha\ln(|\mathcal{P}|)})$ if assumption (A0) is satisfied by the observation graphs $G_t^O, \forall t \in [T]$.

The choices of the parameters β and η (which are non-trivial in the cases where the observation graphs are non-symmetric) that yield these results will be given in Appendix F. We also note that a trivial upper-bound of α_t is the number of vertices of the graph G_t^O which is E (the number of edges in G). In general, the more connected G_t^O is, the smaller α may be chosen; and thus the better upper-bound of the expected regret. In the (classical) semi-bandit setting, $\alpha_t = E, \forall t \in [T]$ and in the full-information setting, $\alpha_t = 1, \forall t \in [T]$. Finally, we also note that, if $|\mathcal{P}| = \mathcal{O}(\exp(n))$ (this is typical in practice, including the CB and HS games), the bound in Corollary 3.3-(i) matches in order with the bounds (ignoring the logarithmic factors) given by the FPL-IX algorithm (see [23]). On the other hand, the form of the regret bound provided by the EXP3-IX algorithm (see [23]) does not allow us to compare directly with the bound of EXP3-OE in the general SOPPP. In [23], EXP3-IX is only analyzed when $n = 1$, i.e., $|\mathcal{P}| = E$; in this case, we observe that the bound given by our EXP3-OE algorithm is better than that of EXP3-IX (by some multiplicative constants).

4 Colonel Blotto Games and Hide-and-Seek Games as SOPPP

Given the regret analysis of EXP3-OE in SOPPP, we now return to our main motivation, the Colonel Blotto and the Hide-and-Seek games, and discuss how to apply our findings to these games. To address this, we define formally the online version of the games and show how these problems can be formulated as SOPPP in Sections 4.1 and 4.2, then we demonstrate the benefit of using the EXP3-OE algorithm for learning in these games (Section 4.3).

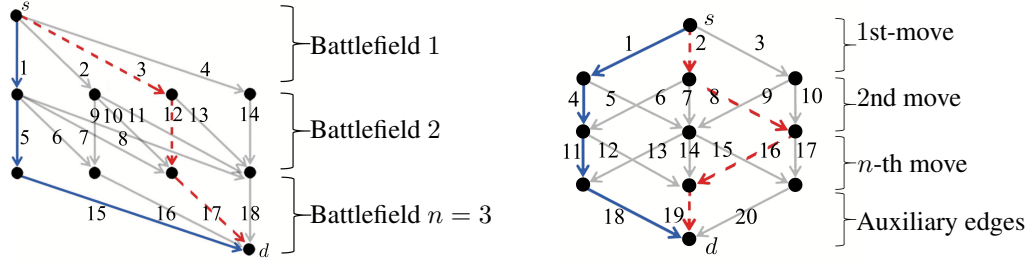
4.1 Colonel Blotto Games as an SOPPP

The online Colonel Blotto game. This is a game between a learner and an adversary over $n \geq 2$ battlefields within a time horizon $T > 0$. Each battlefield $i \in [n]$ has a value $b_t(i) > 0$ (unknown to the learner) at stage t such that $\sum_{i=1}^n b_t(i) = 1$. At stage t , the learner needs to distribute k troops ($k \geq 1$ is fixed) towards the battlefields while the adversary simultaneously allocate hers. The learner's strategy set is $S_{k,n} := \{\mathbf{x} \in \mathbb{N}^n : \sum_{i=1}^n x(i) = k\}$. At stage t and battlefield $i \in [n]$, if the adversary's allocation is strictly larger than the learner's allocation, the learner loses this battlefield and she suffers the loss $b_t(i)$; if they have tie allocations, she suffers the loss $b_t(i)/2$; otherwise, she wins and suffers no loss. At the end of stage t , the learner observes the loss from each battlefield (and which battlefield she wins, ties, or loses) but not the adversary's allocations. The learner's loss at each time is the sum of the losses from all the battlefields. The objective of the learner is then to minimize her loss over a finite period of time.

While this problem can be formulated as a standard OCOMB, it is difficult to derive an efficient learning algorithm under that formulation, due to the learner's exponentially large set of strategies that she can choose from per stage. Instead, we show that by reformulating the problem as an SOPPP, we will be able to exploit the advantages of the EXP3-OE algorithm to solve it. To do so, first note that the learner can deduce several side-observations as follows: (i) if she allocates $x_t(i)$ troops to battlefield i and wins, she knows that if she had allocated more than $x_t(i)$ troops to i , she would also have won; (ii) if she knows the allocations are tie at battlefield i , she knows exactly the adversary's allocation to this battlefield and deduce all the losses she might have suffered if she had allocated differently to battlefield i ; (iii) if she allocates $x_t(i)$ troops to battlefield i and loses, she knows that if she had allocated less than $x_t(i)$ to battlefield i , she would also have lost.

Now, to cast the CB game as SOPPP, for each instance of the parameters k and n , we create a DAG $G := G_{k,n}$ such that the strategy set $S_{k,n}$ has a one-to-one correspondence to the paths set \mathcal{P} of $G_{k,n}$. The formal definition of $G_{k,n}$ will be given in Appendix G; due to the lack of space, we only present here an example illustrating the graph of an instance of the CB game in Figure 1-(a). The graph $G_{k,n}$ has $E = \mathcal{O}(k^2 n)$ edges and $|\mathcal{P}| = |S_{k,n}| = \Omega(2^{\min\{n-1, k\}})$ paths while the length of every path is n . Each edge in $G_{k,n}$ corresponds to allocating a certain amount of troops to a battlefield. Therefore, the CB game model is equivalent to a PPP where at each stage the learner chooses a path in $G_{k,n}$ and the loss on each edge is generated from the allocations of the adversary and the learner (corresponding to that edge) according to the rules of the game. At stage t , the (semi-bandit) feedback and the side-observations⁶ deduced by the learner as described above infers an observation graph G_t^O . This formulation indeed transforms any CB game into an SOPPP.

Note that since there are edges in $G_{m,n}$ that refer to the same allocation (e.g., the edges 5, 9, 12, and 14 in $G_{3,3}$ all refer to allocating 0 troops to battlefield 2), in the observation graphs, the vertices corresponding to these edges are always connected. Therefore, an upper bound of the independence number α_t of G_t^O in the CB game is $\alpha_{CB} = n(k+1) = \mathcal{O}(nk)$. Moreover, we can verify that the observation graph G_t^O of the CB game satisfies assumption (A0) for any t and it is *non-symmetric*.



(a) The graph $G_{3,3}$ corresponding to the CB game with $k = n = 3$. E.g., the bold-blue path represents the strategy $(0, 0, 3)$ while the dash-red path represents the strategy $(2, 0, 1)$.

(b) The graph $G_{3,3,1}$ corresponding to the HS game with $k = n = 3$ and $\kappa = 1$. E.g., the blue-bold path represents the $(1, 1, 1)$ search and the red-dashed path represents the $(2, 3, 2)$ search.

Figure 1: Examples of the graphs corresponding to the CB game and the HS game.

4.2 Hide-and-Seek Games as an SOPPP

The online Hide-and-Seek game. This is a repeated game (within the time horizon $T > 0$) between a hider and a seeker. In this work, we consider that the learner plays the role of the seeker and the hider is the adversary. There are k locations, indexed from 1 to k . At each stage t , the learner sequentially chooses n locations, called an n -search, to seek for the hider, that is, she chooses an $\mathbf{x}_t \in [k]^n$ (if $\mathbf{x}_t(i) = j$, we say that location j is her i -th move). The hider maliciously assigns losses on all k locations (intuitively, these losses can be the wasted time supervising a mismatch location or the probability that the hider does not hide there, etc.). In this work, we consider the following condition on how the hider/adversary assigns the losses on the locations.

(C1) At stage t , the adversary secretly assigns a loss $\mathbf{b}_t(j)$ to each location $j \in [k]$ (unknown to the learner). These losses are fixed throughout the n -search of the learner.

The learner's loss at stage t is the sum of the losses from her chosen locations in the n -search at stage t , that is $\sum_{i \in [n], j \in [k]} \mathbb{I}_{\{\mathbf{x}_t(i)=j\}} \mathbf{b}_t(j)$. Moreover, often in practice the n -search of the learner needs to satisfy some constraints. In this work, as an example, we use the following constraint: $|\mathbf{x}_t(i) - \mathbf{x}_t(i+1)| \leq \kappa, \forall i \in [n]$ for a fixed $\kappa \in [0, k-1]$ (called the *coherence constraint*), i.e., the seeker cannot search too far away from her previously chosen location.⁷ At the end of stage t , the

⁶E.g., in Figure 1-(a), if the learner chooses a path going through edge 10 (corresponding to allocating 1 troop to battlefield 2) and wins (thus, the loss at edge 10 is 0), then she deduces that the losses on the edges 6, 7, 8, 10, 11, and 13 (corresponding to allocating at least 1 troop to battlefield 2) are all 0.

⁷Our results can be applied to HS games with other constraints, such as $\mathbf{x}_t(i) \leq \mathbf{x}_t(i+1), \forall i \in [n]$, i.e., she can only search forward; or, $\sum_{i \in [n]} \mathbb{I}_{\{\mathbf{x}_t(i)=k^*\}} \leq \kappa$, i.e., she cannot search a location $k^* \in [k]$ more than κ times, etc.

learner only observes the losses from the locations she chose among her n -search, and her objective is to minimize her total loss over T .

Similar to the case of the CB game, tackling the HS game as a standard OCOMB is computationally involved. As such, we follow the SOPPP formulation instead. In particular, knowing that the adversary follows condition (C1), the learner can deduce the following side-observations: within a stage, the loss at each location remains the same no matter when it is chosen among the n -search; that is, knowing the loss of choosing location j as her i -th move, the learner knows all the loss if she chooses location j as her i' -th move for any $i' \neq i$. Given this, we create a DAG $G := G_{k,n,\kappa}$ whose paths set has a one-to-one correspondence to the set containing all feasible n -search of the learner in the HS game with k locations under κ -coherent constraint. A formal definition of $G_{k,n,\kappa}$ is given in Appendix G. The HS game is equivalent to the PPP where the learner chooses a path in $G_{k,n,\kappa}$ and edges' losses are generated by the adversary at each stage (note that to ensure all paths end at d , there are n auxiliary edges in $G_{k,n,\kappa}$ that are always embedded with 0 losses). Figure 1-(b) illustrates the corresponding graph of an instance of the HS game. We note that there are $E = \mathcal{O}(k^2 n)$ edges and $|\mathcal{P}| = \Omega(\kappa^{n-1})$ paths in $G_{k,n,\kappa}$.

The semi-bandit feedback and side-observations as described above generate an observation graph G_t^O at time t (e.g., in Figure 1-(b), the edges 1, 4, 6, 11, and 13 represent that location 1 is chosen; thus, they mutually reveal each other). The independence number of G_t^O is $\alpha_{\text{HS}} = k$ for any t . We note that the observation graphs of the HS game are *symmetric* and *do not satisfy* assumption (A0). Finally, we consider a relaxation of condition (C1):

(C2) *At stage t , the adversary assigns a loss $b_t(j)$ on each location $j \in [k]$. For $i = 2, \dots, n$, after the learner chooses, say location j_i , as her i -th move, the adversary can observe that and change the losses $b_t(j)$ for any location that has not been searched before by the learner;⁸ i.e., she can change the losses $b_t(j), \forall j \notin \{j_1, \dots, j_i\}$.*

By replacing condition (C1) with condition (C2), we can limit the side-observations of the learner: she can only deduce that if $i_1 < i_2$, the edges in $G_{k,n,\kappa}$ representing choosing a location as the i_1 -th move reveals the edges representing choosing that same location as the i_2 -th move; but *not vice versa*. In this case, the observation graph G_t^O only contains *directed* edges; however, its independence number is still $\alpha_{\text{HS}} = k$ as in the HS games with condition (C1).

4.3 Performance of EXP3-OE in the Colonel Blotto and Hide-and-Seek Games

Having formulated the CB game and the HS game as SOPPPs, we can use the EXP3-OE algorithm to achieve the following results (deduced directly from Corollary 3.3).

Corollary 4.1. *The expected regret of the EXP3-OE algorithm satisfies:*

- (i) $R_T \leq \tilde{\mathcal{O}}(\sqrt{nT\alpha_{\text{CB}} \ln(|\mathcal{P}|)}) = \tilde{\mathcal{O}}(\sqrt{Tn^3k})$ in the CB games with k troops and n battlefields.
- (ii) $R_T \leq \tilde{\mathcal{O}}(n\sqrt{T\alpha_{\text{HS}} \ln(|\mathcal{P}|)}) = \tilde{\mathcal{O}}(\sqrt{Tn^3k})$ in the HS games with k locations and n -search.

At a high-level, given the same scale on their inputs, the independence numbers of the observation graphs in HS games are smaller than in CB games (by a multiplicative factor of n). However, since assumption (A0) is satisfied by the observation graphs of the CB games and not by the HS games, the expected regret bounds of the EXP3-OE algorithm in these games have the same order of magnitude. From Corollary 4.1, we note that in the CB games, the order of the regret bounds given by EXP3-OE is better than that of the FLP-IX algorithm (thanks to the fact that (A0) is satisfied). On the other hand, in the HS games with condition (C1) involving symmetric observation graphs, the regret bounds of the EXP3-OE algorithm improves the bound of FLP-IX but they are still in the same order of the games' parameters (ignoring the logarithmic factors). Finally, we compare the regret guarantees given by our EXP3-OE algorithm and by the Online Stochastic Mirror Descent algorithm (henceforth, OSMD; see [3])—the benchmark algorithm for OCOMB with semi-bandit feedback (although OSMD does not run efficiently in general). Applying OSMD to the CB and HS games (as SOPPP), the side-observations are ignored and the expected regret bound guaranteed by OSMD is in $\mathcal{O}(\sqrt{TnE}) = \mathcal{O}(\sqrt{Tn^2k^2})$. Using the parameters β and η chosen for Corollary 3.3 and 4.1 (see Appendix F) in the corresponding cases of the observation graphs, the EXP3-OE algorithm provides a better upper-bound of the expected regret than OSMD in the CB games if $\mathcal{O}\left(n \cdot \ln(n^3k^5\sqrt{T})\right) \leq k$;

⁸An interpretation is that by searching a location, the learner/seeker “discovers and secures” that location; therefore, the adversary/hider cannot change her assigned loss at that place.

in the HS games with condition (C1) if $\mathcal{O}(n \ln \kappa) \leq k$; and in the HS games with condition (C2) if $n \cdot \ln \kappa \ln(n^4 k^5 \sqrt{T}) \leq \mathcal{O}(k)$. A proof of this statement is given in Appendix H.

5 Conclusion

In this work, we introduce the EXP3-OE algorithm for the path planning problem with semi-bandit feedback and side-observations. EXP3-OE is always efficiently implementable. Moreover, it matches the regret guarantees compared to that of the FPL-IX algorithm. We apply our findings to derive the first solutions to the online version of the Colonel Blotto and Hide-and-Seek games. This work also extends the scope of application of the PPP model in practice, even for large instances.

References

- [1] Noga Alon, Nicolo Cesa-Bianchi, Ofer Dekel, and Tomer Koren. Online learning with feedback graphs: Beyond bandits. In *JMLR Workshop and Conference Proceedings*, volume 40. Microtome Publishing, 2015.
- [2] Noga Alon, Nicolo Cesa-Bianchi, Claudio Gentile, and Yishay Mansour. From bandits to experts: A tale of domination and independence. In *Advances in Neural Information Processing Systems*, pages 1610–1618, 2013.
- [3] Jean-Yves Audibert, Sébastien Bubeck, and Gábor Lugosi. Regret in online combinatorial optimization. *Mathematics of Operations Research*, 39(1):31–45, 2014.
- [4] Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. Gambling in a rigged casino: The adversarial multi-armed bandit problem. In *focs*, page 322. IEEE, 1995.
- [5] Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. The nonstochastic multiarmed bandit problem. *SIAM journal on computing*, 32(1):48–77, 2002.
- [6] Soheil Behnezhad, Sina Dehghani, Mahsa Derakhshan, MohammadTaghi HajiAghayi, and Saeed Seddighin. Faster and simpler algorithm for optimal strategies of Blotto game. In *AAAI*, pages 369–375, 2017.
- [7] Lilian Besson and Emilie Kaufmann. What doubling tricks can and can’t do for multi-armed bandits. *arXiv preprint arXiv:1803.06971*, 2018.
- [8] Sourabh Bhattacharya, Tamer Başar, and Maurizio Falcone. Surveillance for security as a pursuit-evasion game. In *International Conference on Decision and Game Theory for Security*, pages 370–379. Springer, 2014.
- [9] Sourabh Bhattacharya and Seth Hutchinson. On the existence of nash equilibrium for a two player pursuit-evasion game with visibility constraints. In *Algorithmic Foundation of Robotics VIII*, pages 251–265. Springer, 2009.
- [10] Jeremiah Blocki, Nicolas Christin, Anupam Datta, Ariel D Procaccia, and Arunesh Sinha. Audit games. In *Twenty-Third International Joint Conference on Artificial Intelligence*, 2013.
- [11] Emile Borel. La théorie du jeu et les équations intégrales à noyau symétrique. *Comptes rendus de l’Académie des Sciences*, 173(1304-1308):58, 1921.
- [12] Joseph L Bower and Clark G Gilbert. *From resource allocation to strategy*. Oxford University Press, 2005.
- [13] Nicolo Cesa-Bianchi and Gábor Lugosi. Combinatorial bandits. *Journal of Computer and System Sciences*, 78(5):1404–1422, 2012.
- [14] Wei Chen, Yajun Wang, and Yang Yuan. Combinatorial multi-armed bandit: General framework and applications. In *International Conference on Machine Learning*, pages 151–159, 2013.
- [15] Pern Hui Chia. Colonel Blotto in web security. In *The Eleventh Workshop on Economics and Information Security, WEIS Rump Session*, pages 141–150, 2012.
- [16] Timothy H Chung, Geoffrey A Hollinger, and Volkan Isler. Search and pursuit-evasion in mobile robotics. *Autonomous robots*, 31(4):299, 2011.
- [17] Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences*, 55(1):119–139, 1997.
- [18] Oliver Gross and Robert Wagner. A continuous Colonel Blotto game. U.S.Air Force Project RAND Research Memorandum, 1950.

- [19] JD Grote. *The theory and application of differential games*. Springer, 1975.
- [20] András György, Tamás Linder, Gábor Lugosi, and György Ottucsák. The on-line shortest path problem under partial monitoring. *Journal of Machine Learning Research*, 8(Oct):2369–2403, 2007.
- [21] Joao P Hespanha, Maria Prandini, and Shankar Sastry. Probabilistic pursuit-evasion games: A one-step nash approach. In *Proceedings of the 39th IEEE Conference on Decision and Control (Cat. No. 00CH37187)*, volume 3, pages 2272–2277. IEEE, 2000.
- [22] Adam Kalai and Santosh Vempala. Efficient algorithms for online decision problems. *Journal of Computer and System Sciences*, 71(3):291–307, 2005.
- [23] Tomáš Kocák, Gergely Neu, Michal Valko, and Rémi Munos. Efficient learning by implicit exploration in bandit problems with side observations. In *Advances in Neural Information Processing Systems*, pages 613–621, 2014.
- [24] Dmytro Korzhyk, Vincent Conitzer, and Ronald Parr. Complexity of computing optimal stackelberg strategies in security resource allocation games. In *Twenty-Fourth AAAI Conference on Artificial Intelligence*, 2010.
- [25] Dan Kovenock and Brian Roberson. Coalitional Colonel Blotto games with application to the economics of alliances. *Journal of Public Economic Theory*, 14(4):653–676, 2012.
- [26] Shie Mannor and Ohad Shamir. From bandits to experts: On the value of side-observations. In *Advances in Neural Information Processing Systems*, pages 684–692, 2011.
- [27] Antonia Maria Masucci and Alonso Silva. Strategic resource allocation for competitive influence in social networks. In *Allerton*, pages 951–958, 2014.
- [28] Antonia Maria Masucci and Alonso Silva. Defensive resource allocation in social networks. In *CDC*, pages 2927–2932, 2015.
- [29] Vishnu Navda, Aniruddha Bohra, Samrat Ganguly, and Dan Rubenstein. Using channel hopping to increase 802.11 resilience to jamming attacks. In *INFOCOM 2007. 26th IEEE International Conference on Computer Communications*. IEEE, pages 2526–2530. IEEE, 2007.
- [30] Brian Roberson. The Colonel Blotto game. *Economic Theory*, 29(1):2–24, 2006.
- [31] Shinsaku Sakaue, Masakazu Ishihata, and Shin-ichi Minato. Efficient bandit combinatorial optimization algorithm with zero-suppressed binary decision diagrams. In *International Conference on Artificial Intelligence and Statistics*, pages 585–594, 2018.
- [32] Galina Schwartz, Patrick Loiseau, and Shankar S Sastry. The heterogeneous Colonel Blotto game. In *NetGCoop*, pages 232–238, 2014.
- [33] Eiji Takimoto and Manfred K Warmuth. Path kernels and multiplicative updates. *Journal of Machine Learning Research*, 4(Oct):773–818, 2003.
- [34] Rene Vidal, Omid Shakernia, H Jin Kim, David Hyunchul Shim, and Shankar Sastry. Probabilistic pursuit-evasion games: theory, implementation, and experimental evaluation. *IEEE transactions on robotics and automation*, 18(5):662–669, 2002.
- [35] John Von Neumann. A certain zero-sum two-person game equivalent to the optimal assignment problem. *Contributions to the Theory of Games*, 2:5–12, 1953.
- [36] Dong Quan Vu, Patrick Loiseau, and Alonso Silva. Efficient computation of approximate equilibria in discrete Colonel Blotto games. In *IJCAI-ECAI*, July 2018.
- [37] Qingsi Wang and Mingyan Liu. Learning in hide-and-seek. *IEEE/ACM Transactions on Networking*, 24(2):1279–1292, 2016.
- [38] Wenyuan Xu, Wade Trappe, Yanyong Zhang, and Timothy Wood. The feasibility of launching and detecting jamming attacks in wireless networks. In *Proceedings of the 6th ACM international symposium on Mobile ad hoc networking and computing*, pages 46–57. ACM, 2005.
- [39] Yaakov Yavin. Pursuit–evasion differential games with deception or interrupted observation. In *Pursuit-Evasion Differential Games*, pages 191–203. Elsevier, 1987.
- [40] Chongjie Zhang, Victor Lesser, and Prashant Shenoy. A multi-agent learning approach to online distributed resource allocation. In *Twenty-First International Joint Conference on Artificial Intelligence*, 2009.

A Weight Pushing for Path Sampling

We re-visit some useful results in the literature. In this section, we consider a DAG G with parameters as introduced in Section 2. For simplicity, we assume that each edge in \mathcal{E} belongs to at least one path in \mathcal{P} . Let us respectively denote by $C(u)$ and $F(u)$ the set of the direct successors and the set of the direct predecessors of any vertex $u \in \mathcal{V}$. Moreover, let $e_{[u,v]}$ and $\mathcal{P}_{u,v}$ respectively denote the edge and the set of all paths from vertex u to vertex v . Let us consider a weight $w(e) > 0$ for each edge $e \in \mathcal{E}$. It is needed in the EXP3-OE algorithm to sample a path $\tilde{\mathbf{p}} \in \mathcal{P}$ with the probability:

$$d(\tilde{\mathbf{p}}) := \left[\prod_{e \in \tilde{\mathbf{p}}} w(e) \right] / \left[\sum_{\mathbf{p} \in \mathcal{P}} \prod_{e \in \mathbf{p}} w(e) \right]. \quad (5)$$

A direct computation and sampling from $d_t(\tilde{\mathbf{p}})$, $\forall \tilde{\mathbf{p}} \in \mathcal{P}$ takes $\mathcal{O}(|\mathcal{P}|)$ time which is very inefficient. To efficiently sample the path, we first label the vertices set by $\mathcal{V} = \{s = u_0, u_1, \dots, d = u_K\}$ such that if there exists an edge connecting u_i to u_j then $i < j$. We then define the following terms for each vertex $u \in \mathcal{V}$:

$$H(s, u) := \sum_{\mathbf{p} \in \mathcal{P}_{s,u}} \prod_{e \in \mathbf{p}} w(e) \text{ and } H(u, d) := \sum_{\mathbf{p} \in \mathcal{P}_{u,d}} \prod_{e \in \mathbf{p}} w(e).$$

Intuitively, $H(u, v)$ is the aggregate weight of all paths from vertex u to vertex v and $H(s, d)$ is exactly the denominator in (5). These terms $H(s, u)$ and $H(u, d)$, $\forall u \in \mathcal{V}$ can be recursively computed by Algorithm 3 that runs in $\mathcal{O}(E)$ time, through dynamic programming. This technique is called *weight pushing* and can be found in [20, 31, 33].

Algorithm 3 Weight Pushing.

- 1: **Input:** Graph G , set of weights $\{w(e), e \in \mathcal{E}\}$.
 - 2: Initialization $H(s, u_0) := H(u_K, d) := 1$.
 - 3: **for** $k = 1$ **to** K **do**
 - 4: $H(u_{K-k}, d) := \sum_{v \in C(u_{K-k})} w(e_{[u_{K-k}, v]}) H(v, d)$.
 - 5: $H(s, u_k) := \sum_{v \in F(u_k)} w(e_{[v, u_k]}) H(s, v)$.
 - 6: **end for**
 - 7: **Output:** $H(s, u), H(u, d), \forall u \in \mathcal{V}$.
-

Based on Algorithm 3, we construct Algorithm 4 that uses the weights $w(e), e \in \mathcal{E}$ as inputs and randomly outputs a path in \mathcal{P} . Intuitively, starting from the root vertex $s = u_0$, Algorithm 4 sequentially samples vertices by vertices based on the terms $H(u, v)$ computed by Algorithm 3. It is noteworthy that Algorithm 4 also runs in $\mathcal{O}(E)$ time and it is trivial to prove that the probability that a path \mathbf{p} is sampled from Algorithm 4 matches exactly $d(\mathbf{p})$.

Algorithm 4 Path-sampling Algorithm.

- 1: **Input:** Graph G , set of weights $\{w(e), e \in \mathcal{E}\}$.
 - 2: $H(u, d), \forall u \in \mathcal{V}$ are computed by Algorithm 3.
 - 3: Initialize $\mathbf{Q} := \{s\}$, vertex $u := s$.
 - 4: **while** $u \neq d$ **do**
 - 5: Sample a vertex v from $C(u)$ with probability $w(e_{[u,v]}) H(v, d) / H(u, d)$.
 - 6: Add v to the set \mathbf{Q} and update $u := v$.
 - 7: **end while**
 - 8: **Output:** $\tilde{\mathbf{p}} \in \mathcal{P}$ going through all the vertices in \mathbf{Q}
-

B Proof of Algorithm 2's Output

Proof. Fixing an edge $e \in \mathcal{E}$, we prove that when Algorithm 2 takes the edges weights $\{w_t(e), e \in \mathcal{E}\}$ as the input, it outputs exactly $q_t = \sum_{\mathbf{p} \in \mathbb{O}_t(e)} d_t(\mathbf{p})$. We note that if $e' \in \mathfrak{R}_t(e) := \{e' : e' \rightarrow e\}$, then $\{\mathbf{p} \in \mathcal{P} : \mathbf{p} \ni e'\} \subset \mathbb{O}_t(e)$.

We denote $|\mathfrak{R}_t(e)| = \rho_e$ and label the edges in the set $\mathfrak{R}_t(e)$ by $\{e_1, e_2, \dots, e_{\rho_e}\}$. We let the for-loop in lines 3–8 of Algorithm 2 consecutively run with the edges in $R_t(e)$ as follows:

(i) After the for-loop runs for e_1 , we have $K(e_1) := \sum_{\mathbf{p} \ni e_1} \prod_{\bar{e} \in \mathbf{p}} \bar{w}(\bar{e}) = \sum_{\mathbf{p} \ni e_1} w_t(\mathbf{p})$; therefore, $q_t(e) = \sum_{\mathbf{p} \ni e_1} d_t(\mathbf{p})$ since $H(s, d) = \sum_{\mathbf{p} \in \mathcal{P}} w_t(\mathbf{p})$ computed from the original weights $w_t(\bar{e}), \bar{e} \in \mathcal{E}$. Due to line 8 that sets $\bar{w}(e_1) := 0$, henceforth in Algorithm 2, the weight $\bar{w}(\mathbf{p}) := \prod_{e \in \mathbf{p}} \bar{w}(e)$ of any path \mathbf{p} that contains e_1 is set to 0.

(ii) Let the for-loop run for e_2 , we have $K(e_2) := \sum_{\mathbf{p} \ni e_2} \bar{w}(\mathbf{p}) = \sum_{\{\mathbf{p} \ni e_2\} \setminus \{\mathbf{p} \ni e_1\}} w_t(\mathbf{p})$ because any path $\mathbf{p} \ni e_1$ has the weight $\bar{w}(\mathbf{p}) = 0$. Therefore, $q_t(e) = \sum_{\mathbf{p} \ni e_1} d_t(\mathbf{p}) + \sum_{\{\mathbf{p} \ni e_2\} \setminus \{\mathbf{p} \ni e_1\}} d_t(\mathbf{p})$.

(iii) Similarly, after the for-loop runs for e_i (where $i \in \{3, \dots, \rho_e\}$), we have:

$$q_t(e) = \sum_{k=1}^i \left(\sum_{\{\mathbf{p} \ni e_k\} \setminus \bigcup_{j < k} \{\mathbf{p} \ni e_j\}} d_t(\mathbf{p}) \right).$$

(iv) Therefore, after the for-loop finishes running for every edge in $\mathfrak{R}_t(e)$; we have $q_t := \sum_{\mathbf{p} \in \mathcal{O}_t(e)} d_t(\mathbf{p})$ where each term $d_t(\mathbf{p})$ was only counted once even if \mathbf{p} contains more than one edge that reveals the edge e .

□

C Proof of Theorem 3.1

Theorem 3.1. *The expected regret of the EXP3-OE algorithm in the SOPPP satisfies:*

$$R_T \leq \ln(|\mathcal{P}|)/\eta + [\beta + (n \cdot \eta)/2] \cdot \sum_{t \in [T]} Q_t. \quad (4)$$

Proof. We first denote⁹ $W_t := \sum_{\mathbf{p} \in \mathcal{P}} w_t(\mathbf{p}), \forall t \in [T]$. From line 9 of Algorithm 1, we trivially have:

$$w_{t+1}(\mathbf{p}) = w_t(\mathbf{p}) \cdot \exp(-\eta \hat{L}_t(\mathbf{p})), \forall \mathbf{p} \in \mathcal{P}, \forall t \in [T-1]. \quad (6)$$

Here, we recall $\hat{L}_t(\mathbf{p}) := \sum_{e \in \mathbf{p}} \hat{\ell}_t(e)$, then from (2), we have:

$$\mathbb{E}_t [\hat{L}_t(\mathbf{p})] \leq L_t(\mathbf{p}) := \sum_{e \in \mathbf{p}} \ell_t(e), \forall \mathbf{p} \in \mathcal{P}. \quad (7)$$

Under the condition that $0 < \eta$, we obtain:

$$\begin{aligned} \frac{W_{t+1}}{W_t} &= \sum_{\mathbf{p} \in \mathcal{P}} \frac{w_{t+1}(\mathbf{p})}{W_t} = \sum_{\mathbf{p} \in \mathcal{P}} \frac{w_t(\mathbf{p}) \cdot \exp(-\eta \hat{L}_t(\mathbf{p}))}{W_t} \\ &= \sum_{\mathbf{p} \in \mathcal{P}} d_t(\mathbf{p}) \cdot \exp(-\eta \hat{L}_t(\mathbf{p})) \\ &\leq \sum_{\mathbf{p} \in \mathcal{P}} \left[d_t(\mathbf{p}) \left(1 - \eta \hat{L}_t(\mathbf{p}) + \frac{\eta^2}{2} (\hat{L}_t(\mathbf{p}))^2 \right) \right] \\ &= 1 - \sum_{\mathbf{p} \in \mathcal{P}} \left[d_t(\mathbf{p}) \left(\eta \hat{L}_t(\mathbf{p}) - \frac{\eta^2}{2} (\hat{L}_t(\mathbf{p}))^2 \right) \right]. \end{aligned} \quad (8)$$

Here, the second equality comes from (6) and the inequality comes from the fact that $\exp(-x) \leq 1 - x + x^2/2$ for $x := \eta \hat{L}_t(\mathbf{p}) \geq 0$. From (8) and the inequality $\ln(1 - y) \leq -y$ for any $y < 1$, we have the following inequality:¹⁰

$$\ln \left(\frac{W_{T+1}}{W_1} \right) = \sum_{t=1}^T \ln \left(\frac{W_{t+1}}{W_t} \right) \leq \sum_{t=1}^T \left(-\eta \sum_{\mathbf{p} \in \mathcal{P}} d_t(\mathbf{p}) \hat{L}_t(\mathbf{p}) + \frac{\eta^2}{2} \sum_{\mathbf{p} \in \mathcal{P}} d_t(\mathbf{p}) (\hat{L}_t(\mathbf{p}))^2 \right). \quad (9)$$

⁹We recall that $w_t(\mathbf{p}) := \prod_{e \in \mathbf{p}} w_t(e)$.

¹⁰We can easily check that $\eta \hat{L}_t(\mathbf{p}) - \eta^2 \hat{L}_t(\mathbf{p})^2 / 2 < 1$ for any $\eta > 0$.

On the other hand, let us fix a path $\mathbf{p}^* \in \mathcal{P}$, then

$$\begin{aligned} \ln \left(\frac{W_{T+1}}{W_1} \right) &\geq \ln \left(\frac{w_{T+1}(\mathbf{p}^*)}{W_1} \right) = \ln \frac{w_T(\mathbf{p}^*) \exp(-\eta \hat{L}_T(\mathbf{p}^*))}{|\mathcal{P}|} \\ &= \ln \frac{w_{T-1}(\mathbf{p}^*) \exp(-\eta \hat{L}_T(\mathbf{p}^*) - \eta \hat{L}_{T-1}(\mathbf{p}^*))}{|\mathcal{P}|} \\ &= -\eta \sum_{t=1}^T \hat{L}_t(\mathbf{p}^*) - \ln(|\mathcal{P}|). \end{aligned} \quad (10)$$

In the arguments leading to (10), we again use (6) and the fact that $w_1(\mathbf{p}) = 1, \forall \mathbf{p} \in \mathcal{P}$, including $w_1(\mathbf{p}^*)$. Therefore, combining (9) and (10) then dividing both sides by η , we have that

$$\sum_{t=1}^T \sum_{\mathbf{p} \in \mathcal{P}} d_t(\mathbf{p}) \hat{L}_t(\mathbf{p}) \leq \frac{\ln(|\mathcal{P}|)}{\eta} + \sum_{t=1}^T \hat{L}_t(\mathbf{p}^*) + \frac{\eta}{2} \sum_{t=1}^T \sum_{\mathbf{p} \in \mathcal{P}} d_t(\mathbf{p}) (\hat{L}_t(\mathbf{p}))^2. \quad (11)$$

Now, we take the expectation \mathbb{E}_t w.r.t. to the randomness in choosing $\tilde{\mathbf{p}}_t$ on (11), then we apply (7) to obtain:

$$\sum_{t=1}^T \sum_{\mathbf{p} \in \mathcal{P}} d_t(\mathbf{p}) \mathbb{E}_t[\hat{L}_t(\mathbf{p})] \leq \frac{\ln(|\mathcal{P}|)}{\eta} + \sum_{t=1}^T L_t(\mathbf{p}^*) + \frac{\eta}{2} \sum_{t=1}^T \sum_{\mathbf{p} \in \mathcal{P}} d_t(\mathbf{p}) \mathbb{E}_t[\hat{L}_t(\mathbf{p})^2]. \quad (12)$$

Now, we look for a lower bound of $\sum_{\mathbf{p} \in \mathcal{P}} d_t(\mathbf{p}) \mathbb{E}_t[\hat{L}_t(\mathbf{p})]$. For any fixed $\mathbf{p} \in \mathcal{P}$, we consider:

$$\begin{aligned} \mathbb{E}_t \left[\sum_{e \in \mathbf{p}} \hat{\ell}_t(e) \right] &= \sum_{\tilde{\mathbf{p}} \in \mathcal{P}} \left[d_t(\tilde{\mathbf{p}}) \sum_{e \in \mathbf{p}} \left(\frac{\ell_t(e)}{q_t(e) + \beta} \mathbb{I}_{\{e \in \mathbb{O}_t(\tilde{\mathbf{p}})\}} \right) \right] \\ &= \sum_{e \in \mathbf{p}} \sum_{\tilde{\mathbf{p}} \in \mathbb{O}(e)} d_t(\tilde{\mathbf{p}}) \frac{\ell_t(e)}{q_t(e) + \beta} \\ &= \sum_{e \in \mathbf{p}} \frac{q_t(e) \ell_t(e)}{q_t(e) + \beta}. \end{aligned} \quad (13)$$

On the other hand, applying (13) and recalling that $\ell_t(e) \leq 1, \forall e \in \mathcal{E}$, we have:

$$\begin{aligned} \sum_{\mathbf{p} \in \mathcal{P}} d_t(\mathbf{p}) \mathbb{E}_t[\hat{L}_t(\mathbf{p})] - \sum_{\mathbf{p} \in \mathcal{P}} d_t(\mathbf{p}) L_t(\mathbf{p}) &= \sum_{\mathbf{p} \in \mathcal{P}} d_t(\mathbf{p}) \sum_{e \in \mathbf{p}} \frac{q_t(e) \ell_t(e)}{q_t(e) + \beta} - \sum_{\mathbf{p} \in \mathcal{P}} d_t(\mathbf{p}) \sum_{e \in \mathbf{p}} \ell_t(e) \\ &= \sum_{\mathbf{p} \in \mathcal{P}} d_t(\mathbf{p}) \sum_{e \in \mathbf{p}} \ell_t(e) \left(\frac{q_t(e)}{q_t(e) + \beta} - 1 \right) \\ &\geq - \sum_{\mathbf{p} \in \mathcal{P}} d_t(\mathbf{p}) \sum_{e \in \mathbf{p}} \frac{\beta}{q_t(e) + \beta} \\ &= - \beta \sum_{e \in \mathcal{E}} \frac{\sum_{\mathbf{p} \ni e} d_t(\mathbf{p})}{q_t(e) + \beta} \\ &= - \beta Q_t. \end{aligned} \quad (14)$$

Now, we look for an upper bound of $\sum_{\mathbf{p} \in \mathcal{P}} d_t(\mathbf{p}) \mathbb{E}_t[\hat{L}_t(\mathbf{p})^2]$. To do this, we fix $\mathbf{p} \in \mathcal{P}$ and consider

$$\begin{aligned} \mathbb{E}_t[\hat{L}_t(\mathbf{p})^2] &= \mathbb{E}_t \left[\left(\sum_{e \in \mathbf{p}} \hat{\ell}_t(e) \right)^2 \right] \leq n \cdot \mathbb{E}_t \left[\sum_{e \in \mathbf{p}} \hat{\ell}_t(e)^2 \right] \\ &= n \cdot \sum_{\tilde{\mathbf{p}} \in \mathcal{P}} \left[d_t(\tilde{\mathbf{p}}) \sum_{e \in \mathbf{p}} \left(\frac{\ell_t(e)}{q_t(e) + \beta} \mathbb{I}_{\{e \in \mathbb{O}_t(\tilde{\mathbf{p}})\}} \right)^2 \right] \end{aligned}$$

$$\begin{aligned}
&\leq n \cdot \sum_{e \in \mathcal{P}} \sum_{\tilde{\mathbf{p}} \in \mathbb{O}_t(e)} d_t(\tilde{\mathbf{p}}) \frac{1}{(q_t(e) + \beta)^2} \\
&= n \cdot \sum_{e \in \mathcal{P}} q_t(e) \frac{1}{(q_t(e) + \beta)^2} \\
&\leq n \cdot \sum_{e \in \mathcal{P}} \frac{1}{q_t(e) + \beta}.
\end{aligned} \tag{15}$$

The first inequality comes from applying Cauchy–Schwarz inequality. The second inequality comes from the fact that $\ell_t(e) \leq 1$ and the last inequality comes from $q_t(e) \leq q_t(e) + \beta$ since $\beta > 0$.

Now, applying (15), we can bound

$$\begin{aligned}
\sum_{\mathbf{p} \in \mathcal{P}} d_t(\mathbf{p}) \mathbb{E}_t \left[\hat{L}_t(\mathbf{p})^2 \right] &\leq n \cdot \sum_{\mathbf{p} \in \mathcal{P}} d_t(\mathbf{p}) \sum_{e \in \mathcal{P}} \frac{1}{q_t(e) + \beta} \\
&= n \cdot \sum_{e \in \mathcal{E}} \sum_{\mathbf{p} \ni e} d_t(\mathbf{p}) \frac{1}{q_t(e) + \beta} \\
&= n \cdot \sum_{e \in \mathcal{E}} \frac{r_t(e)}{q_t(e) + \beta} = n \cdot Q_t.
\end{aligned} \tag{16}$$

Here, we recall the notation $r_t(e)$ and Q_t defined in (3). Replacing (14) and (16) into (12), we have that the following inequality holds for any $\mathbf{p}^* \in \mathcal{P}$.

$$\begin{aligned}
&\sum_{t=1}^T \sum_{\mathbf{p} \in \mathcal{P}} d_t(\mathbf{p}) L_t(\mathbf{p}) - \sum_{t=1}^T \beta Q_t - \sum_{t=1}^T L_t(\mathbf{p}^*) \leq \frac{\ln(|\mathcal{P}|)}{\eta} + \frac{\eta}{2} \sum_{t=1}^T n Q_t \\
&\Rightarrow R_T \leq \frac{\ln(|\mathcal{P}|)}{\eta} + \sum_{t=1}^T Q_t \left(n \frac{\eta}{2} + \beta \right).
\end{aligned}$$

□

D Lemmas on Graphs' Independence Numbers

In this section, we present some lemmas in graph theory that will be used in the next section to prove Theorem 3.2. Consider a graph \tilde{G} whose vertices set and edges set are respectively denoted by $\tilde{\mathcal{V}}$ and $\tilde{\mathcal{E}}$. Let $\tilde{\alpha}$ be its independence number.

Lemma D.1. *Let \tilde{G} be an directed graph and I_v be the in-degree of the vertex $v \in \tilde{\mathcal{V}}$, then*

$$\sum_{v \in \tilde{\mathcal{V}}} [1/(1 + I_v)] \leq 2\tilde{\alpha} \ln \left(1 + |\tilde{\mathcal{V}}|/\tilde{\alpha} \right).$$

A proof of this lemma can be found in Lemma 10 of [2].

Lemma D.2. *Let \tilde{G} be a directed graph with self-loops and consider the numbers $k(v) \in [0, 1], \forall v \in \tilde{\mathcal{V}}$ such that there exists $\gamma > 0$ and $\sum_{v \in \tilde{\mathcal{V}}} k(v) \leq \gamma$. For any $c > 0$, we have*

$$\sum_{v \in \tilde{\mathcal{V}}} \frac{k(v)}{\frac{1}{\gamma} \sum_{v' \rightarrow v} k(v') + c} \leq 2\gamma\tilde{\alpha} \ln \left(1 + \frac{\gamma[|\tilde{\mathcal{V}}|^2/c] + |\tilde{\mathcal{V}}|}{\tilde{\alpha}} \right) + 2\gamma.$$

A proof of this lemma can be found in Lemma 1 of [23].

Lemma D.3. *Let \tilde{G} be an undirected graph with self-loops and consider the numbers $k(v) \geq 0, v \in \tilde{\mathcal{V}}$. We have*

$$\sum_{v \in \tilde{\mathcal{V}}} \left[k(v) / \sum_{v' \rightarrow v} k(v') \right] \leq \tilde{\alpha}.$$

This lemma is extracted from Lemma 3 of [26].

E Proof of Theorem 3.2

Theorem 3.2. Let us denote $M := \lceil 2E^2/\beta \rceil$, $N_t := \ln\left(1 + \frac{M+E}{\alpha_t}\right)$ and $K_t := \ln\left(1 + \frac{nM+E}{\alpha_t}\right)$. Upper-bounds of Q_t in different cases of G_t^O are given in the following table:

	SATISFIES (A0)	NOT SATISFIES (A0)
SYMMETRIC	α_t	$n\alpha_t$
NON-SYMMETRIC	$1+2\alpha_t N_t$	$2n(1+\alpha_t K_t)$

Case 1: G_t^O does not satisfy assumption (A0). Fixing an edge e , due to the fact that n is the length of the longest paths in \mathcal{P} , we have

$$nq_t(e) = n \sum_{\mathbf{p} \in \mathbb{O}_t(e)} d_t(\mathbf{p}) \geq \sum_{e' \rightarrow e} \sum_{\mathbf{p} \ni e'} d_t(\mathbf{p}) = \sum_{e' \rightarrow e} r_t(e') \quad (17)$$

$$\Rightarrow Q_t = \sum_{e \in \mathcal{E}} \frac{r_t(e)}{q_t(e) + \beta} \leq \sum_{e \in \mathcal{E}} \frac{r_t(e)}{\frac{1}{n} \sum_{e' \rightarrow e} r_t(e') + \beta}. \quad (18)$$

Case 1.1: If G_t^O is a non-symmetric (i.e., directed) graph, we apply Lemma D.2 with $\gamma = n, c = \beta$ on the graph $\tilde{G} = G_t^O$ (whose vertices set $\tilde{\mathcal{V}}$ corresponds to the edges set \mathcal{E} of G) and the numbers¹¹ $k(v_e) = r_t(e), \forall v_e \in \tilde{\mathcal{V}}$ (i.e., $\forall e \in \mathcal{E}$). We obtain the following inequality:

$$\sum_{e \in \mathcal{E}} \frac{r_t(e)}{\frac{1}{n} \sum_{e' \rightarrow e} r_t(e') + \beta} \leq 2n\alpha_t \ln\left(1 + \frac{n\lceil E^2/\beta \rceil + E}{\alpha_t}\right) + 2n.$$

Case 1.2: If G_t^O is a symmetric (i.e. undirected) graph, we apply Lemma D.3 with the graph $\tilde{G} = G_t^O$ (whose vertices set $\tilde{\mathcal{V}}$ corresponds to the edges set \mathcal{E} of the graph G) and the numbers $k(v_e) = r_t(e), \forall v_e \in \tilde{\mathcal{V}}$ (i.e., $\forall e \in \mathcal{E}$) to obtain:

$$\sum_{e \in \mathcal{E}} \frac{r_t(e)}{\frac{1}{n} \sum_{e' \rightarrow e} r_t(e') + \beta} \leq n \sum_{e \in \mathcal{E}} \frac{r_t(e)}{\sum_{e' \rightarrow e} r_t(e')} \leq n\alpha_t.$$

Case 2: G_t^O satisfies assumption (A0). Under this assumption, $q_t(e) = \sum_{e' \rightarrow e} r_t(e')$ due to the definition of $\mathbb{O}_t(e)$. Therefore, $Q_t = \sum_{e \in \mathcal{E}} [r_t(e) / (\sum_{e' \rightarrow e} r_t(e') + \beta)]$.

Case 2.1: If G_t^O is a non-symmetric (i.e., directed) graph. We consider a discretized version of $d_t(\mathbf{p})$ for any path $\mathbf{p} \in \mathcal{P}$ that is $\tilde{d}_t(\mathbf{p}) := k/M$ where k is the unique integer such that $(k-1)/M \leq d_t(\mathbf{p}) \leq k/M$; thus, $\tilde{d}_t(\mathbf{p}) - 1/M \leq d_t(\mathbf{p}) \leq \tilde{d}_t(\mathbf{p})$.

Let us denote the discretized version of $r_t(e)$ by $\tilde{r}_t(e) := \sum_{\mathbf{p} \ni e} \tilde{d}_t(\mathbf{p})$. We deduce that $r_t(e) \leq \tilde{r}_t(e)$ and

$$\sum_{e' \rightarrow e} r_t(e) \geq \sum_{e' \rightarrow e} \left(\tilde{r}_t(e') - \frac{1}{M} \right) \geq \sum_{e' \rightarrow e} \tilde{r}_t(e') - \frac{E}{M}.$$

We obtain the bound:

$$Q_t = \sum_{e \in \mathcal{E}} \frac{r_t(e)}{\left(\sum_{e' \rightarrow e} r_t(e') + \beta \right)} \leq \sum_{e \in \mathcal{E}} \frac{\tilde{r}_t(e)}{\sum_{e' \rightarrow e} \tilde{r}_t(e') - E/M + \beta}. \quad (19)$$

¹¹We verify that these numbers satisfy

$$\sum_{e \in \mathcal{E}} r_t(e) = \sum_{e \in \mathcal{E}} \sum_{\mathbf{p} \ni e} d_t(\mathbf{p}) = \sum_{\mathbf{p} \in \mathcal{P}} \sum_{e \in \mathbf{p}} d_t(\mathbf{p}) \leq \sum_{\mathbf{p} \in \mathcal{P}} n d_t(\mathbf{p}) = n.$$

We now consider the following inequality: If $a, b \geq 0$ and $a + b \geq B > A > 0$, then

$$\frac{a}{a+b-A} \leq \frac{a}{a+b} + \frac{A}{B-A}. \quad (20)$$

A proof of this inequality can be found in Lemma 12 of [2]. Applying (20)¹² with $a = \tilde{r}_t(e)$, $b = \sum_{e' \rightarrow e, e' \neq e} \tilde{r}_t(e') + \beta$, $A = E/M$, and $B = \beta$ to (19), we have

$$Q_t \leq \sum_{e \in \mathcal{E}} \left(\frac{\tilde{r}_t(e)}{\sum_{e' \rightarrow e} \tilde{r}_t(e') + \beta} + \frac{E/M}{\beta - E/M} \right) \leq \sum_{e \in \mathcal{E}} \frac{\tilde{r}_t(e)}{\sum_{e' \rightarrow e} \tilde{r}_t(e')} + 1. \quad (21)$$

The last inequality comes from the fact that $\frac{E}{M\beta-E} \leq \frac{E}{2E^2-E} \leq \frac{1}{2E-1} \leq \frac{1}{E}$, $\forall E \geq 1$.

Finally, we create an auxiliary graph G_t^* such that:

- (i) Corresponding to each edge e in G (i.e., each vertex v_e in G_t^O), there is a clique, called $\mathbb{C}(e)$, in the auxiliary graph G_t^* with $M\tilde{r}_t(e) \in \mathbb{N}$ vertices.
- (ii) In each clique $\mathbb{C}(e)$ of G_t^* , all vertices are pairwise connected with length-two cycles. That is, for any $k, k' \in \mathbb{C}(e)$, there is an edge from k to k' and there is an edge from k' to k in G_t^* .
- (iii) If $e \rightarrow e'$, i.e., there is an edge in G_t^O connecting v_e and $v_{e'}$; then in G_t^* , all vertices in the clique $\mathbb{C}(e)$ are connected to all vertices in $\mathbb{C}(e')$.

We observe that the independence number α_t of G_t^O is equal to the independence number of G_t^* . Moreover, the in-degree of each vertex $k \in (e)$ in the graph G_t^* is:

$$I_k^* = M\tilde{r}_t(e) - 1 + \sum_{e' \rightarrow e, e' \neq e} M\tilde{r}_t(e') = \sum_{e' \rightarrow e} M\tilde{r}_t(e') - 1. \quad (22)$$

Let us denote V_t^* the set of all vertices in G_t^* , then we have:

$$\begin{aligned} \sum_{e \in \mathcal{E}} \frac{\tilde{r}_t(e)}{\sum_{e' \rightarrow e} \tilde{r}_t(e')} &= \sum_{e \in \mathcal{E}} \frac{M\tilde{r}_t(e)}{\sum_{e' \rightarrow e} M\tilde{r}_t(e')} = \sum_{e \in \mathcal{E}} \sum_{k \in \mathbb{C}(e)} \frac{1}{I_k^* + 1} \\ &= \sum_{k \in V_t^*} \frac{1}{I_k^* + 1} \leq 2\alpha_t \ln \left(1 + \frac{M+E}{\alpha_t} \right). \end{aligned} \quad (23)$$

Here, the second equality comes from the fact that $|\mathbb{C}(e)| = M\tilde{r}_t(e)$ and (22). The inequality is obtained by applying Lemma D.1 to the graph G_t^* and the fact that $|V_t^*| = \sum_{e \in \mathcal{E}} M\tilde{r}_t(e) \leq M \sum_{e \in \mathcal{E}} (r_t(e) + 1/M) \leq E + M$.

In conclusion, combining (21) and (23), we obtain the regret-upper bound as given in Theorem 3.2 for this case of the observation graph.

Case 2.2: Finally, if G_t^O is a symmetric (i.e., undirected) graph, we again apply Lemma D.3 to the graph $\tilde{G} = G_t^O$ and the numbers $k(v_e) = r_t(e)$ to obtain that $Q_t \leq \sum_{e \in \mathcal{E}} [r_t(e) / \sum_{e' \rightarrow e} r_t(e')] \leq \alpha_t$. \square

F Parameters Tuning for EXP3-OE: Proof of Corollary 3.3

In this section, we suggest a choice of β and η that guarantees the expected regret given in Corollary 3.3.

Corollary 3.3. *In SOPPP, let α be an upper bound of $\alpha_t, \forall t \in [T]$. With appropriate choices of the parameters η and β , the expected regret of the EXP3-OE algorithm is:*

- (i) $R_T \leq \tilde{O}(n\sqrt{T\alpha \ln(|\mathcal{P}|)})$ in the general cases.
- (ii) $R_T \leq \tilde{O}(\sqrt{nT\alpha \ln(|\mathcal{P}|)})$ if assumption (A0) is satisfied by the observation graphs $G_t^O, \forall t \in [T]$.

¹²Trivially, we can verify that $a + b \geq B$ and $B > A$ comes from the fact that $\beta \geq \beta \frac{1}{E} > \frac{E}{\lceil 2E^2/\beta \rceil}$.

Case 1: Non-symmetric (i.e. directed) observation graphs that do not satisfy assumption (A0). We find the parameters β and η such that $R_t \leq \tilde{\mathcal{O}}\left(n\sqrt{T\alpha}\right)$. We note that $\alpha_t \geq 1, \forall t \in [T]$; therefore, recalling that α is an upper bound of α_t , from Theorem 3.1 and 3.2, we have:

$$\begin{aligned} R_T &\leq \frac{\ln(|\mathcal{P}|)}{\eta} + \sum_{t=1}^T \left(n\frac{\eta}{2} + \beta\right) 2n \left[1 + \alpha_t \ln\left(1 + \frac{nM+E}{\alpha_t}\right)\right] \\ &\leq \frac{\ln(|\mathcal{P}|)}{\eta} + T \left(n\frac{\eta}{2} + \beta\right) 2n [1 + \alpha \ln(\alpha + nM + E)] \\ &= \frac{\ln(|\mathcal{P}|)}{\eta} + \eta T n^2 [1 + \alpha \ln(\alpha + nM + E)] + 2\beta T n [1 + \alpha \ln(\alpha + nM + E)]. \end{aligned} \quad (24)$$

Recalling that $M := \lceil 2E^2/\beta \rceil$, by choosing any

$$\beta \leq 1/\sqrt{Tn[1 + \alpha \ln(\alpha + n\lceil E^2/\beta \rceil + E)]}, \quad (25)$$

$$\text{and } \eta = \sqrt{\ln |\mathcal{P}|} / \sqrt{n^2 T [1 + \alpha \ln(\alpha + n\lceil E^2/\beta \rceil + E)]},$$

we obtain the bound:

$$\begin{aligned} R_T &\leq 2n\sqrt{T \ln |\mathcal{P}| \cdot [1 + \alpha \ln(\alpha + nM + E)]} + 2\sqrt{Tn[\alpha + \alpha \ln(\alpha + nM + E)]} \\ &\leq \tilde{\mathcal{O}}\left(n\sqrt{T\alpha \ln(|\mathcal{P}|)}\right). \end{aligned} \quad (26)$$

In practice, as long as it satisfies (25), the larger β is, the better upper-bounds that EXP3-OE gives.

Finally, as an example that (25) always has at least one solution, we now prove that it holds with

$$\beta^* = \frac{-Tn^2 E^2 + \sqrt{(Tn^2 E^2)^2 + 4Tn(1 + \alpha \ln \alpha + E + n)}}{2Tn(1 + \alpha \ln \alpha + E + n)}. \quad (27)$$

Indeed, $\beta^* > 0$ and it satisfies:

$$\begin{aligned} &\beta^{*2} \cdot Tn(1 + \alpha \ln \alpha + E + n) + \beta^* Tn^2 E^2 = 1. \\ \Rightarrow &\beta^{*2} \cdot Tn(1 + \alpha \ln \alpha + E) + \beta^{*2} Tn^2 \left(\frac{E^2}{\beta^*} + 1\right) = 1 \\ \Rightarrow &\beta^{*2} \cdot Tn(1 + \alpha \ln \alpha + E) + \beta^{*2} Tn^2 \left[\frac{E^2}{\beta^*}\right] \leq 1 \\ \Rightarrow &\beta^* \leq \frac{1}{\sqrt{Tn(1 + \alpha \ln \alpha + E + nM)}}. \end{aligned}$$

On the other hand, applying the inequality $\ln(1+x) \leq x, \forall x \geq 0$, we have:

$$\begin{aligned} &\frac{nM + E}{\alpha} \geq \ln\left(1 + \frac{nM + E}{\alpha}\right) \\ \Rightarrow &\frac{nM + E}{\alpha} + \ln \alpha \geq \ln(\alpha + nM + E) \\ \Rightarrow &nM + E + \alpha \ln \alpha + 1 \geq \alpha \ln(\alpha + nM + E) + 1 \\ \Rightarrow &\frac{1}{\sqrt{Tn(1 + \alpha \ln \alpha + nM + E)}} \leq \frac{1}{\sqrt{Tn(\alpha \ln(\alpha + nM + E) + 1)}}. \end{aligned}$$

Therefore, β^* satisfies (25).

Case 2: non-symmetric observation graphs G_t^O satisfying assumption (A0), $\forall t$. We will prove that $R_T \leq \tilde{\mathcal{O}}\left(\sqrt{nT\alpha \ln(|\mathcal{P}|)}\right)$ for any

$$\beta \leq 1/\sqrt{T\alpha[1 + 2 \ln(1 + \lceil E^2/\beta \rceil + E)]}, \quad (28)$$

$$\eta = 2\sqrt{\ln |\mathcal{P}|} / \sqrt{Tn\alpha[1 + 2 \ln(\alpha + M + E)]}. \quad (29)$$

Indeed, from Theorem 3.1 and 3.2, we have:

$$\begin{aligned}
R_T &\leq \frac{\ln(|\mathcal{P}|)}{\eta} + \sum_{t=1}^T \left(n \frac{\eta}{2} + \beta \right) \left[1 + 2\alpha_t \ln \left(1 + \frac{M+E}{\alpha_t} \right) \right] \\
&\leq \frac{\ln(|\mathcal{P}|)}{\eta} + \sum_{t=1}^T \left(n \frac{\eta}{2} + \beta \right) [\alpha + 2\alpha \ln(1 + M + E)] \\
&= \frac{\ln(|\mathcal{P}|)}{\eta} + \eta T \alpha \frac{n}{2} [1 + 2 \ln(1 + M + E)] + \beta T \alpha [1 + 2 \ln(1 + M + E)]. \quad (30)
\end{aligned}$$

We replace (28) and (29) into (30) and obtain:

$$\begin{aligned}
R_T &\leq \frac{3}{2} \sqrt{T n \alpha [1 + 2 \ln(1 + M + E)] \cdot \ln |\mathcal{P}|} + \sqrt{T \alpha [1 + 2 \ln(1 + M + E)]}. \quad (31) \\
&\leq \tilde{\mathcal{O}} \left(\sqrt{n \alpha T \ln(|\mathcal{P}|)} \right).
\end{aligned}$$

A choice for β that satisfies (28) is

$$\beta^* := \frac{-T \alpha E^2 + \sqrt{(T \alpha E^2)^2 + T \alpha (3 + 2E)}}{T \alpha (3 + 2E)}. \quad (32)$$

Case 3: symmetric observation graphs that do not satisfy (A0). Trivially, we have that if $\beta := 1/\sqrt{n \alpha T}$ and $\eta = 2\sqrt{\ln |\mathcal{P}|/\sqrt{n^2 \alpha T}}$, then

$$\begin{aligned}
R_T &\leq \frac{\ln |\mathcal{P}|}{\eta} + \left(n \frac{\eta}{2} + \beta \right) n \alpha T \\
&= \frac{1}{2} n \sqrt{\alpha T \ln |\mathcal{P}|} + n \sqrt{\alpha T \ln |\mathcal{P}|} + \sqrt{n \alpha T} \\
&\leq \tilde{\mathcal{O}} \left(n \sqrt{\alpha T \ln(|\mathcal{P}|)} \right). \quad (33)
\end{aligned}$$

Case 4: all observation graphs are symmetric and satisfy (A0). From Theorem 3.1 and 3.2, we trivially have that if $\beta := 1/\sqrt{\alpha T}$ and $\eta = 2\sqrt{\ln |\mathcal{P}|/\sqrt{n \alpha T}}$, then $R_T \leq \tilde{\mathcal{O}} \left(\sqrt{n \alpha T \ln(|\mathcal{P}|)} \right)$.

G Graphical Representation of the Games' Actions Sets

G.1 The Actions Set of the Colonel Blotto Games

We give a description of the graph corresponding to the actions set of the learner in the CB game who distributes k troops to n battlefields.

Definition G.1 (CB Graph). *The graph $G_{k,n}$ is a DAG that contains:*

- (i) $N := 2 + (k+1)(n-1)$ vertices arranged into $n+1$ layers. Layer 0 and Layer n , each contains only one vertex, respectively labeled $s := (0,0)$ —the source vertex and $d := (n,k)$ —the destination vertex. Each Layer $i \in [n-1]$ contains $k+1$ vertices whose labels are ordered from left to right by $(i,0), (i,1), \dots, (i,k)$.
- (ii) There are directed edges from vertex $(0,0)$ to every vertex in Layer 1 and edges from every vertex in Layer $n-1$ to vertex (n,k) . For $i \in \{1, 2, \dots, n-2\}$, there exists an edge connecting vertex (i, j_1) (of Layer i) to vertex $(i+1, j_2)$ (of Layer $(i+1)$) if $k \geq j_2 \geq j_1 \geq 0$.

Particularly, $G_{k,n}$ has $E = (k+1)[4+(n-2)(k+2)]/2 = \mathcal{O}(nk^2)$ edges and $|\mathcal{P}| = \binom{n+k-1}{n-1} = \mathcal{O}(2^{\min\{n-1, k\}})$ paths going from vertex $s := (0,0)$ to vertex $d := (n,k)$. The edge connecting vertex (i, j_1) to vertex $(i+1, j_2)$ for any $i \in \{0, 1, \dots, n-1\}$ represents allocating $(j_2 - j_1)$ troops to battlefield $i+1$. Moreover, each path from s to d represents a strategy in $S_{k,n}$. This is formally stated in Proposition G.2.

Proposition G.2. *Given k and n , there is a one-to-one mapping between the action set $S_{k,n}$ of the learner in the CB game (with k troops and n battlefields) and the set of all paths from vertex s to vertex d of the graph $G_{k,n}$.*

The proof of this proposition is trivial and can be intuitively seen in Figure 1-(a). We note that a similar graph can be found in [6]; however, it is used for a completely different purpose and it also contains more edges and paths than $G_{k,n}$ (that are not useful in this work).

G.2 The Actions Set of the Hide-and-Seek game

We give a description of the graph corresponding to the actions set of the learner in the HS games with the n -search among k locations and coherence constraints $|\mathbf{x}_t(i) - \mathbf{x}_t(i+1)| \leq \kappa, \forall i \in [n]$ for a fixed $\kappa \in [0, k-1]$.

Definition G.3 (HS Graph). *The graph $G_{k,\kappa,n}$ is a DAG that contains:*

- (i) $N := 2 + kn$ vertices arranged into $n + 2$ layers. Layer 0 and Layer $(n + 1)$, each contains only one vertex, respectively labeled s —the source vertex and d —the destination vertex. Each Layer $i \in \{1, \dots, n\}$ contains k vertices whose labels are ordered from left to right by $(i, 1), (i, 2), \dots, (i, k)$.
- (ii) There are directed edges from vertex s to every vertex in Layer 1 and edges from every vertex in Layer n to vertex d . For $i \in \{1, 2, \dots, n-1\}$, there exists an edge connecting vertex (i, j_1) to vertex $(i+1, j_2)$ if $|j_1 - j_2| \leq \kappa$.

The graph $G_{k,\kappa,n}$ has $E = 2k + (n-1)[k + \kappa(2k - \kappa - 1)] = \mathcal{O}(nk^2)$ edges and at least $\Omega(\kappa^{n-1})$ paths from s to d . The edges ending at vertex d are the auxiliary edges that are added just to guarantee that all paths end at d ; these edges do not represent any intuitive quantity related to the game. For the remaining edges, any edge that ends at the vertex (i, j) represents choosing the location j as the i -th move. In other words, a path starting from s , passing by vertices $(1, j_1), (2, j_2), \dots, (n, j_n)$ and ending at d represents the n -search that chooses location j_1 , then moves to location j_2 , then moves to location j_3 , and so on.

Proposition G.4. *Given k, κ and n , there is a one-to-one mapping between the action set $S_{k,\kappa,n}$ of the learner in the HS game (with n -search among k locations and coherence constraints with parameter κ) and the set of all paths from vertex s to vertex d of the graph $G_{k,\kappa,n}$.*

H EXP3-OE Algorithm versus OSMD Algorithm in the CB and HS Games

(i) As stated in Section 4, the observation graphs in the CB games are non-symmetric and they satisfy assumption (A0). If we choose $\beta = \beta^*$ as in (32), then β satisfies (28). Moreover, $\beta = \mathcal{O}(1/\sqrt{TnE})$; thus, $M = \mathcal{O}(E^2\sqrt{TnE})$. From (31), the expected regret of EXP3-OE in this case is bounded by $\mathcal{O}(\sqrt{Tn(\alpha_{CB} \ln M \ln |\mathcal{P}|)})$ (recall that $\alpha_{CB} = kn$ is an upper bound of independence numbers of the observation graphs in the CB games). Therefore, to guarantee that this bound is better than the bound of the OSMD algorithm (that is $\sqrt{2TnE}$), the following inequality needs to hold:

$$\begin{aligned}
& \mathcal{O}(\alpha_{CB} \cdot \ln M \ln |\mathcal{P}|) \leq E \\
& \Rightarrow \mathcal{O}\left(nk \cdot \ln(E^2\sqrt{TnE}) \ln(2^n)\right) \leq nk^2 \\
& \Rightarrow \mathcal{O}\left(\ln(E^2\sqrt{TnE}) \ln(2^n)\right) \leq k \\
& \Rightarrow \mathcal{O}\left(n \ln(n^3 k^5 \sqrt{T})\right) \leq k.
\end{aligned}$$

(ii) As stated in Section 4, the observation graphs in the HS games with condition (C1) are symmetric and do not satisfy assumption (A0). If we choose $\beta = 1/\sqrt{n\alpha T}$ then by (33), we have that R_T is bounded by $\mathcal{O}\left(n\sqrt{\alpha_{HS}T \ln |\mathcal{P}|}\right)$ (recall that $\alpha_{HS} = k$ is an upper bound of the independence numbers of the observation graphs in the HS games). Therefore, to guarantee that this bound is better than the bound of the OSMD algorithm in HS games, the following inequality needs to hold:

$$\begin{aligned}
& \mathcal{O}(\alpha_{HS} \cdot n \ln |\mathcal{P}|) \leq E \\
& \Rightarrow \mathcal{O}(k \cdot n \ln |\mathcal{P}|) \leq nk^2 \\
& \Rightarrow \mathcal{O}(\ln |\mathcal{P}|) \leq k \\
& \Rightarrow \mathcal{O}(n \ln \kappa) \leq k.
\end{aligned}$$

(iii) Finally, the observation graphs in the HS games with condition (C2) are non-symmetric and do not satisfy assumption (A0). Therefore, if we choose $\beta = \beta^*$ as in (27), then β satisfies (25). In this case, $\beta = \mathcal{O}(1/\sqrt{TnE})$ and $M = \mathcal{O}(E^2\sqrt{TnE})$. Therefore, from (26), in this case, R_T is bounded by $\mathcal{O}(n\sqrt{T\alpha_{HS}\ln\alpha_{HS}\ln(nM)})$. Therefore, to guarantee that this bound is better than the bound of OSMD (that is, $\sqrt{2TnE}$), the following inequality needs to hold:

$$\begin{aligned} \mathcal{O}(\alpha_{HS} \cdot n \ln nM \ln |\mathcal{P}|) &\leq E \\ \Rightarrow \mathcal{O}\left(nk \ln(\kappa^n) \ln(nE^2\sqrt{TnE})\right) &\leq nk^2 \\ \Rightarrow \mathcal{O}\left(n \ln \kappa \ln(n^4 k^5 \sqrt{T})\right) &\leq k. \end{aligned}$$