



**HAL**  
open science

# The Non-Zero-Sum Game of Steganography in Heterogeneous Environments

Quentin Giboulot, Tomáš Pevný, Andrew D Ker

► **To cite this version:**

Quentin Giboulot, Tomáš Pevný, Andrew D Ker. The Non-Zero-Sum Game of Steganography in Heterogeneous Environments. *IEEE Transactions on Information Forensics and Security*, 2023, 18, pp.4436 - 4448. 10.1109/tifs.2023.3295945 . hal-04463469

**HAL Id: hal-04463469**

**<https://hal.science/hal-04463469v1>**

Submitted on 17 Feb 2024

**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.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

# The Non-Zero-Sum Game of Steganography in Heterogeneous Environments

Quentin Giboulot<sup>1</sup>, Tomáš Pevný, and Andrew D. Ker<sup>2</sup>, *Senior Member, IEEE*

**Abstract**—The highly heterogeneous nature of images found in real-world environments, such as online sharing platforms, has been one of the long-standing obstacles to the transition of steganalysis techniques outside the laboratory. Recent advances in identifying the properties of images relevant to steganalysis as well as the effectiveness of deep neural networks on highly heterogeneous datasets have laid some groundwork for resolving this problem. Despite this progress, we argue that the way the game played between the steganographer and the steganalyst is currently modeled lacks some important features expected in a real-world environment: 1) the steganographer can adapt her cover source choice to the environment and/or to the steganalyst’s classifier, 2) the distribution of cover sources in the environment impacts the optimal threshold for a given classifier, and 3) the steganalyst and steganographer have different goals, hence different utilities. We propose to take these facts into account using a two-player non-zero-sum game constrained by an environment composed of multiple cover sources. We then show how to convert this non-zero-sum game into an equivalent zero-sum game, allowing us to propose two methods to find Nash equilibria for this game: a standard method using the double oracle algorithm and a minimum regret method based on approximating a set of atomistic classifiers. Applying these methods to contemporary steganography and steganalysis in a realistic environment, we show that classifiers which do not adapt to the environment severely underperform when the steganographer is allowed to select into which cover source to embed.

**Index Terms**—Steganalysis, steganography, game theory.

## I. INTRODUCTION

**I**MAGE steganalysis has started its transition from dealing purely with laboratory settings to trying to model the real-world in order to tackle the difficult problem of designing classifiers that can provide performance guarantees in a realistic environment. An illustration is the success of recent steganalysis methodologies during the ALASKA2 competition [1]. Contestants were asked to classify cover and

stego images on a highly heterogeneous dataset containing images coming from many different sources – that is, taken with different cameras, ISO settings and, most importantly, developed with different processing pipelines. The main lesson from this competition was that the use of off-the-shelf state-of-the-art neural network architectures pre-trained on ImageNet and refined on a diverse cover source corpus leads to excellent performance with respect to the metrics that defined the contest [2], [3]. Despite this success, the model under which steganalysis is currently evaluated, even in the case of the ALASKA2 competition, still lacks several features that would be expected when confronting a rational steganographer in a real-world environment.

As a first observation, even though research on cover source mismatch has clearly shown that some cover sources are far more difficult to steganalyze than others [4], [5], the steganographer is never assumed to be able to strategize about the choice of cover source to embed into. In this sense, the steganographer is not considered rational since she will never try to adapt to the environment in order to increase the chances of evasion from the steganalyst. A consequence of this point of view is that the holistic strategy in steganalysis, consisting in training a detector on as many sources as possible, is currently considered the state-of-the-art in steganalysis when dealing with heterogeneous environments, due to its immense success during both ALASKA competitions [1], [6]. One major conclusion of this present work is that such a holistic approach is highly insufficient when facing a rational steganographer.

A realistic model of the game played between the steganographer and the steganalyst should take the environment information into account in order to inform how a rational steganographer should distribute stego objects among the different cover sources. Furthermore, knowledge of the cover source environment is also extremely important to the steganalyst since the proportion of the different sources directly influence the probability of false alarm for a given classifier. For example, a rare cover source will have a very low contribution to the false alarm rate since, irrespective of the strategy of the steganographer, very few innocent users actually use it. As such, the steganalyst should design their classifiers and assign their detection thresholds differently depending on the cover sources distribution in the environment in order to minimize their probability of error. Finally, it should be pointed out that both steganographic algorithms and steganalysis methods are usually evaluated using the same metric. Most often, this metric is chosen to be the minimum probability of error

Manuscript received 14 December 2022; revised 11 April 2023 and 15 June 2023; accepted 9 July 2023. Date of publication 17 July 2023; date of current version 28 July 2023. This work was supported in part by the Operační program Výzkum, vývoj a vzdělávání (OP VVV) Project “Research Center for Informatics” under Grant CZ.02.1.01/0.0/0.0/16\_019/0000765 and in part by the Horizon 2020 European Programme through the UNCOVER Project under Grant 101021687. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Chia-Mu Yu. (Corresponding author: *Quentin Giboulot*.)

Quentin Giboulot and Tomáš Pevný are with the Artificial Intelligence Center, Czech Technical University in Prague, 16000 Prague, Czech Republic (e-mail: gibouloq@protonmail.com).

Andrew D. Ker is with the Department of Computer Science, Oxford University, OX1 3QD Oxford, U.K.

Digital Object Identifier 10.1109/TIFS.2023.3295945

under equal priors,  $P_E$ . Other metrics have been proposed in [7] such as the probability of false alarm at 50% missed detection which is starting to get adoption in the community. However, it is questionable whether both the steganalyst and the steganographer want to minimize the same metric. Indeed, it is not directly obvious why the steganographer should care about the probability of false alarm of the steganalyst since she is only trying to evade detection.

To take these features into account we propose to cast the steganalyst detection problem into a new game theoretic framework. In our setting, the steganalyst tries to minimize both their probability of false alarm and probability of missed detection by selecting from a fixed but possibly infinite set of classifiers. On the other hand, the steganographer has access to a set of cover sources distributed according to a fixed distribution and can choose how to distribute stego objects using these cover sources. In contrast, the steganalyst only wants to maximize the probability of missed detection of the steganalyst and is indifferent to the probability of false alarm.

### A. Related Works

The use of a game theoretic framework in steganography is not new, though its use has been mostly focused on designing better distortion functions for the steganographer. In [8], the authors cast the problem of the steganographer as a zero-sum game played with the steganalyst and show how to formulate a distortion function accordingly in order to reach an equilibrium. In a similar vein, the authors in [9] provide conditions under which adaptivity of the cost function is an optimal strategy in a suitably defined zero-sum game. Finally, and most recently, the Backpack protocol presented in [10], formulates a zero-sum game played between the steganographer and the steganalyst in order to iteratively design a cost function that is optimal for a given dataset of covers. The optimization is performed through the use of an algorithm reminiscent of the classic double oracle algorithm [11] used to iteratively solve minimax optimization problems.

To the best of our knowledge, few prior works in steganalysis have applied game theoretic techniques in order to determine the optimal strategy of the steganalyst for the design of a classifier. We can highlight the insights of [8] and [9] which justified the use of selection-channel aware steganalyzers [12] using game theoretic techniques. On the other hand, there exists some works that address this question outside of steganography such as [13] which studied how a defender should select their classifier thresholds when the attacker can choose to deliver attacks at a given operating point of the ROC curve. Another relevant work is the monograph [14] from Barni and Tondi which fully addressed the question of a game where the attacker selected a function modifying a sequence of objects to evade the defender. In particular they provided a condition in [15, Eq. (3.9)] for the probability distribution of the modified sequence of objects to attain an equilibrium. However, these works are not directly applicable to our problem: in our setting the steganographer does not explicitly control either the probability distributions or the performance of the steganalyst's classifiers. She can only select

which cover sources among a fixed set for embedding which is a more restrictive setting than the one used in these previous works.

The closest work which has tried to integrate the knowledge of the cover source environment into the design of the classifier itself is the recent work of Šepák [16] proposing a formalization of the problem of cover source mismatch. In this work, the authors propose to minimize the maximum regret of a classifier over all cover sources in a given cover source environment. The regret over a source is defined as the difference in  $P_E$  of the classifier with the  $P_E$  of an *atomistic classifier*, that is, a classifier trained exclusively on the source of interest. The min-regret classifier iteratively minimizes the maximum regret by selecting at each optimization step the cover source leading to the highest regret. This implicitly assumes that the steganographer always selects only the worst – in the sense of the hardest to detect with the steganalyst's specific detector – cover source possible. The advantage of this construction is that it does not necessitate knowledge of the proportion of each cover source in the environment. However, we show in Section IV that this definition of the regret cannot guarantee a Nash equilibrium against a rational steganographer.

The ultimate goal of this paper is to leverage our theoretical analysis in order to design steganalysis detectors which can still guarantee strong performance against a rational steganographer.

### B. Contributions and Organizations

In order to address the problems mentioned in the introduction we propose the following contributions, in order of appearance in the paper:

- We formulate a realistic and general game theoretic non-zero-sum model of the game played between the steganographer and the steganalyst taking into account the environment of available cover sources – Section II,
- We then show the equivalence of this non-zero-sum game to a more amenable zero-sum game, thus allowing one to find Nash equilibria using standard techniques – Section III,
- We provide an explicit solution of the zero-sum game in the case where the steganalyst is able to perfectly identify the cover source of all images and has access to an atomistic classifier for each of these cover sources. This result, which is of general interest, also allows us to bound the utility for the more general case where such identification is not possible – Section IV-A,
- We design a fast alternative to standard techniques for finding  $\epsilon$ -Nash equilibria using a minimum regret classifier which tries to match the performance of atomistic classifiers in the perfect identification case – Section IV-C,
- We finally apply the proposed methods to an heterogeneous dataset, using state-of-the-art steganography and steganalysis and compare it to the standard holistic strategy which does not assume a rational steganographer – Section V.

We summarize each of these contributions and their link to each other in Figure 1.

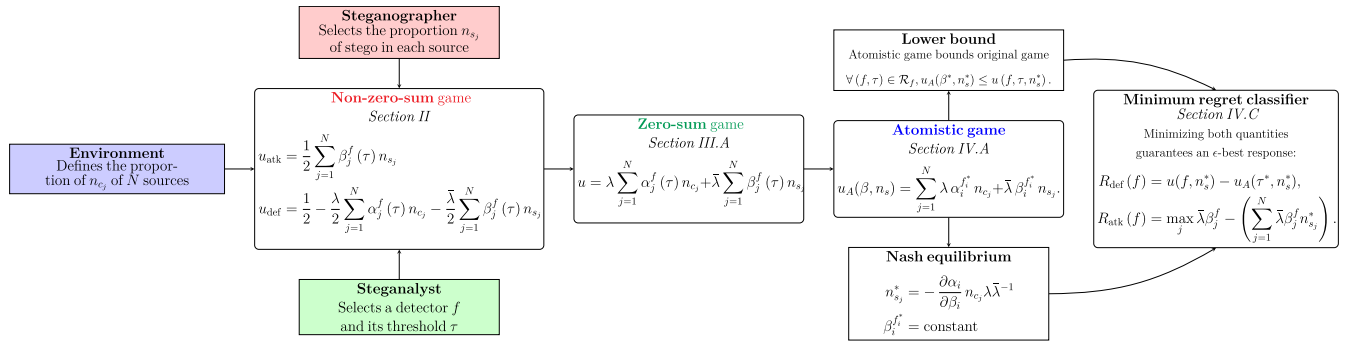


Fig. 1. Summary of the theoretical analysis of this paper as well as its application for the design of classifiers. We start in Section II by formulating the game between the steganographer and the steganalyst as a non-zero-sum game. We then prove in Section III-A that this game is equivalent to a zero-sum game, in the sense that it possesses the same Nash equilibria. We then formulate an approximate version of this game, the atomistic game, in Section IV-A which lower bounds the utility of the zero-sum game at the equilibrium. We provide a Nash equilibrium for this atomistic game and finally, we show in Section IV-C how to leverage this equilibrium to design a classifier which is guaranteed to be an  $\epsilon$ -best response for the steganalyst in the original non-zero-sum game.

## II. GAME FORMULATION

The goal of this section is to provide a formulation of the non-zero-sum game that is played between the steganalyst and the steganographer. We first define the environment in which the game is played as the distribution of cover sources. We then go on to define the strategy set of both the steganographer and steganalyst. Finally, we explain how the game is played and provide a formal definition of the utility function for each player.

### A. The Environment

Outside the laboratory, both the steganographer and the steganalyst have access to a large number of platforms where innocent users share cover objects. The probability distribution of these covers can depend on many parameters. In the case of image steganography, three significant parameters have been identified: the camera and ISO setting used to capture the RAW image and the processing pipeline used to obtain the final cover image [5].

Different combinations of these parameters lead to different probability distributions of covers which we call *cover sources*. We assume that the set  $\mathcal{P}^c$  of cover sources is finite. Even though this might seem unintuitive since the parameters of some image processing algorithms take values in compact subsets of  $\mathbb{R}$ , we follow the recent work of [17] which shows how to find a finite set of representative sources among a possibly infinite set of cover sources.

We finally assume that innocent users draw covers randomly from this finite set of cover sources. However, cover sources are not uniformly distributed in the real-world, hence we follow the work of [16] and model the possible environments as a convex closure of probability distributions in  $\mathcal{P}^c$ :

$$\mathcal{E} \triangleq \left\{ \sum_{i=1}^N n_{c_i} P_i^c \mid \sum_{i=1}^N n_{c_i} = 1, n_{c_i} \geq 0, P_i^c \in \mathcal{P}^c \right\}. \quad (1)$$

At the start of the game, the environment is fixed to one of the member of  $\mathcal{E}$ .

### B. The Steganographer

The goal of the steganographer is to evade the steganalyst as much as possible, in other words, to maximize the probability

of missed detection of the steganalyst. In this paper, we assume a simplified model of the steganographer where she uses a single steganographic algorithm with a single payload for each source. As a consequence, the steganographer has access to  $N$  stego-sources  $P_i^s$ , once again defined as probability distributions. Available strategies consist in choosing the probabilities  $n_{s_i}$  of sampling a stego object from each of these distributions.

Consequently the set of strategies  $\mathcal{R}_s$  of the steganographer is defined by a convex set  $\mathcal{R}_s$  of vectors  $n_s$ :

$$\mathcal{R}_s \triangleq \left\{ (n_{s_i})_{1 \leq i \leq N} \mid \sum_{i=1}^N n_{s_i} = 1, n_{s_i} \geq 0 \right\}. \quad (2)$$

Note that under this formulation, the steganographer is assumed to use mixed strategies since she is expected to select cover sources randomly.

### C. The Steganalyst

Informally, the goal of the steganalyst is to detect as many stego objects as possible while minimizing the probability of false alarms. To do so, they have access to a (possibly infinite) compact set  $\mathcal{K}$  of classifiers. This set can be thought of as a class of classifiers parameterized by a finite number of continuous and bounded parameters.

Let us define the probability of false alarm  $\alpha_j$  and probability of missed detection  $\beta_j$  of a classifier  $f$  for the  $j$ -th source as:

$$\alpha_j^f(\tau) \triangleq \mathbb{E}_{x \sim P_j^c} [f(x) > \tau], \quad (3)$$

$$\beta_j^f(\tau) \triangleq \mathbb{E}_{x \sim P_j^s} [f(x) \leq \tau], \quad (4)$$

where  $\tau \in \bar{\mathbb{R}} = \mathbb{R} \cup \{-\infty, \infty\}$  is a threshold chosen by the steganalyst. The strategy of the steganalyst  $\mathcal{R}_f$  is thus defined as classifier  $f$  and a threshold  $\tau$ :

$$\mathcal{R}_f \triangleq \mathcal{K} \times \bar{\mathbb{R}}. \quad (5)$$

It is important to understand that the steganalyst is only allowed to use a *single classifier* with a single threshold. However, this classifier can itself work in two stages, where the first stage tries to estimate the cover source and a second stage

which uses the best available classifier for the estimated cover source. Nevertheless, the steganalyst does not possess any a priori knowledge about an image cover source, an assumption we relax in Section IV-A.

#### D. The Game

The game is played first with each player choosing their strategies. These choices are fixed until the end of the game.

We assume a repeated game where each round is played by sampling one cover and one stego according to the environment's distribution and the steganographer's strategy. One of these two objects, which we will denote as  $x$ , is then given to the steganalyst at random with equal probabilities. The steganalyst then applies their chosen classifier to this object to decide if it is a cover or a stego. Importantly, the steganalyst only sees a *single image* from an unknown actor at each round. The final utility for each player is the average utility after an infinite number of rounds. We will now define the utility function for each player.

As we have discussed, the steganographer wants to maximize the probability of missed detection of the steganalyst. Note that in this version of the game, the steganographer does not control how many stego images is produced. This leads to the following utility formulation for the steganographer (attacker):

$$\begin{aligned} u_{\text{atk}}(f, \tau, n_s) &\triangleq \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \mathbb{1}[f(x_i) \leq \tau \wedge x_i \sim P^s], \\ &= \frac{1}{2} \mathbb{E}_{x \sim P^s} [f(x) \leq \tau] \\ &= \frac{1}{2} \sum_{j=1}^N \beta_j^f(\tau) n_{s_j}, \end{aligned} \quad (6)$$

where  $P^s$  can be any stego-source.

On the other hand, the steganalyst not only wants to maximize their probability of correctly detecting a stego image, but they also want to minimize their probability of false alarm. Furthermore, they might weight the cost of a false alarm differently than the cost of a missed detection or we might want to model a case where there is an unequal number of cover and stego objects. We reflect this by introducing a weighting parameter  $\lambda \in [0, 1]$  and its complement  $\bar{\lambda} = 1 - \lambda$ . This leads us to the following utility formulation for the steganalyst (defender):

$$\begin{aligned} u_{\text{def}}(f, \tau, n_s) &\triangleq \lim_{n \rightarrow \infty} \frac{\lambda}{n} \sum_{i=1}^n \mathbb{1}[f(x_i) \leq \tau \wedge x_i \sim P^c] \\ &\quad + \lim_{n \rightarrow \infty} \frac{\bar{\lambda}}{n} \sum_{i=1}^n \mathbb{1}[f(x_i) > \tau \wedge x_i \sim P^s] \\ &= \frac{\lambda}{2} \mathbb{E}_{x \sim P^c} [f(x) \leq \tau] + \frac{\bar{\lambda}}{2} \mathbb{E}_{x \sim P^s} [f(x) > \tau] \\ &= \frac{1}{2} - \frac{\lambda}{2} \sum_{j=1}^N \alpha_j^f(\tau) n_{c_j} - \frac{\bar{\lambda}}{2} \sum_{j=1}^N \beta_j^f(\tau) n_{s_j}. \end{aligned} \quad (7)$$

This general formulation of the utility clearly shows that the game played between the steganalyzer and the steganographer is *not a zero-sum game*: the steganographer's utility only depends on the probability of missed detection whereas the steganalyst has a supplementary term depending on the probability of false alarm and each term is weighted differently.

### III. GAME SOLUTIONS

Non-zero-sum games are notoriously difficult to solve. In particular, we loose the property of zero-sum games that minmax solutions, where both players are allowed to randomize their strategies, lead to a Nash equilibrium and that all Nash Equilibria have an identical utility. Thankfully, some non-zero-sum games have a special structure which can be leveraged to convert them to an equivalent zero-sum game with the same equilibria. This is the case for our game, and recasting it as an equivalent zero-sum game allows us, in the rest of paper, to use standard game theoretic techniques to efficiently find its Nash Equilibria.

#### A. Equivalence to a Zero-Sum Game

In our case, the utility of each player can be rewritten as a utility term shared between the two players, denoted as  $u$ , plus a term that is only a function of the strategy of the opposite player. The shared component of the utility can be seen as defining a zero-sum game. By expressing the utility as such, we can show that the best response of each player only depends on the shared component  $u$ . This allows us to conclude that the set of Nash equilibria of the original non-zero-sum game is identical to that of the zero-sum game defined by  $u$ .

Let us now start by defining each component of the non-zero-sum game utility:

$$u'_{\text{def}}(f, \tau) \triangleq \frac{1}{2} - \frac{\lambda}{2} \sum_{j=1}^N \alpha_j^f(\tau) n_{c_j}, \quad (8)$$

$$u(f, \tau, n_s) \triangleq \bar{\lambda} u_{\text{atk}}(f, \tau, n_s) - u'_{\text{def}}(f, \tau). \quad (9)$$

We can then rewrite the utility of both players in Eq. (6-7):

$$u_{\text{atk}}(f, \tau, n_s) = (u(f, \tau, n_s) + u'_{\text{def}}(f, \tau)) \bar{\lambda}^{-1}, \quad (10)$$

$$u_{\text{def}}(f, \tau, n_s) = -u(f, \tau, n_s). \quad (11)$$

Observe that both players share the same utility function  $u$  which depends on the strategy of both players. Furthermore, the other term in the utility of the steganographer only depends on the strategy of steganalyst.

Now let  $(f^*, \tau^*, n_s^*)$  be a Nash equilibrium for the zero-sum game defined by  $u$ . By definition, we then have that  $n_s^*$  is a best response to  $(f^*, \tau^*)$  in the zero-sum game defined by  $u$ . Hence we have that for all  $n_s \in \mathcal{R}_s$ :

$$u(f^*, \tau^*, n_s^*) \geq u(f^*, \tau^*, n_s). \quad (12)$$

Using this inequality, we can upper bound the utility of the steganographer for the original non-zero-sum game  $u_{\text{atk}}$ .

For all  $n_s \in \mathcal{R}_s$ , we have that:

$$\begin{aligned} u_{\text{atk}}(f^*, \tau^*, n_s^*) &= (u(f^*, \tau^*, n_s^*) + u'_{\text{def}}(f^*, \tau^*)) \bar{\lambda}^{-1} \\ &\geq (u(f^*, \tau^*, n_s) + u'_{\text{def}}(f^*, \tau^*)) \bar{\lambda}^{-1} \\ &= u_{\text{atk}}(f^*, \tau^*, n_s). \end{aligned} \quad (13)$$

The key step is the inequality of the second line which is made possible by the fact that the second term of Eq (9),  $u'_{\text{def}}$ , does not depend on  $n_s$  and thus is constant whatever the strategy of the steganographer is.

In the case of the steganalyst, the proof is trivial since the utility  $u_{\text{def}}$  only depends on  $u$ . Thus, for all  $(f, \tau) \in \mathcal{R}_f$ :

$$\begin{aligned} u_{\text{def}}(f^*, \tau^*, n_s^*) &= -u(f^*, \tau^*, n_s^*) \\ &\leq -u(f, \tau, n_s^*) \\ &= u_{\text{def}}(f, \tau, n_s) \end{aligned} \quad (14)$$

Since we have that both  $n_s^*$  and  $(f^*, \tau^*)$  are still best responses in the original non-zero-sum game, we have that  $(f^*, \tau^*, n_s^*)$  is a Nash equilibrium in this non-zero-sum game. Consequently, instead of solving the original game, we can instead find Nash Equilibria in the zero-sum game defined by  $u$ . We emphasize that this is not an approximation: every Nash equilibrium of the zero-sum game defined by  $u$  is equivalent to a Nash equilibrium of the non-zero sum game defined by  $u_{\text{def}}$  and  $u_{\text{atk}}$ , hence computing one equilibrium in one game is equivalent to computing one in the other.

For the rest of the paper, we will drop the  $-1$  constant and the  $\frac{1}{2}$  factor of Eq (9) since they do not affect the equilibrium. The steganography game can then finally be defined by a single utility, denoted as  $u$ , which is expressed as:

$$u(f, \tau, n_s) \triangleq \lambda \sum_{j=1}^N \alpha_j^f(\tau) n_{c_j} + \bar{\lambda} \sum_{j=1}^N \beta_j^f(\tau) n_{s_j}. \quad (15)$$

Notice that  $u$ , which the steganographer aims to maximize and the steganalyst to minimize, equals *the probability of error over the whole environment*. In particular, it is a more general expression of the commonly used  $P_E$ . Note that it is equal to the  $P_E$  when  $\lambda = 0.5$  and when the steganalyst chooses  $\tau$  such that it minimizes the sum of errors.

### B. Double Oracle Solutions for the Zero-Sum Game

A general method to solve large zero-sum games efficiently is to use the double oracle algorithm. This algorithm exists for both finite [11] and continuous [18] games. In the case of finite games, the algorithm is guaranteed to converge to a Nash equilibrium in finitely many iterations [18, Theorem 1.1]. For continuous games, the algorithm is guaranteed to converge to an  $\epsilon$ -equilibrium in finitely many iterations as long as  $\epsilon > 0$  [18, Theorem 1.3] and if we can find the best-response for each player at each iteration. It is important to note that double oracle algorithms *output mixed strategies for both players*. In our case, this means that the steganalyst's strategy under the double oracle solution will not consist in choosing a single classifier but in randomizing the chosen classifier at each round of the game.

In our setting, we cannot rely on computing the complete payoff matrix of the game, due to the possibly infinite size

of the steganalyst strategy set. The use of the double oracle algorithm allows bypassing this limitation by iteratively computing a payoff matrix while also improving upon the strategy of both players at each iteration.

Our game is continuous due to the fact that the steganalyst needs to select a classifier within a compact set and set its threshold. We can thus use the continuous version of the double oracle algorithm as presented in Algorithm 1. The basic idea of the algorithm is to augment, at each iteration, a finite set of pure strategies for both the steganographer and steganalyst. A Nash equilibrium is computed on this finite set using the standard linear programming method. Finally, the best responses for both the steganographer and the steganalyst are found in the original set of strategies  $\mathcal{R}_f$  and  $\mathcal{R}_s$  against this equilibrium. Since *a pure strategy always belongs to the best response set for each player*, we can augment the set of strategies of each player only with pure strategies and repeat until an  $\epsilon$ -Nash equilibrium is reached. From [18, Eq (4), Proposition 2.2], the terminating condition which guarantees an  $\epsilon$ -Nash equilibrium has been reached is that:

$$u(f_j^*, \tau_j^*, (n_s)_{j+1}) - u(f_{j+1}, \tau_{j+1}, (n_s)_j) \leq \epsilon. \quad (16)$$

In words, this means that if, at the  $j$ -th iteration, the best responses of each player –  $(n_s)_{j+1}$  and  $(f_{j+1}, \tau_{j+1})$  – against the equilibrium response of the other player –  $(f_j^*, \tau_j^*)$  and  $(n_s)_j$  – lead to utilities that do not differ more than  $\epsilon$ , then we have reached an  $\epsilon$ -Nash equilibrium.

The only difficulty in the implementation of this algorithm is finding the best responses for both players at each iteration. In the case of the steganographer, this is actually trivial since her number of pure strategies is equal to the number of sources – they are the strategies consisting in setting all  $n_{s_i}$  to zero except for one source. Since the steganographer can restrict the search of best responses to pure strategies, she can simply perform an exhaustive search to find the source leading to the highest probability of missed detection against the current equilibrium strategy of the steganalyst  $(f_j^*, \tau_j^*)$ .

Regarding the steganalyst's best response, it has to be found by finding the classifier which leads to the smallest utility. In practice, this amounts to finding the global minimum for a chosen class of classifier. Since the current state of the art relies on neural networks, for which there is no guarantee to find such a global minimum, the steganalyst will resort to a heuristic that consists of training a classifier on a dataset built using the environment  $n_c$  as well as the steganographer's strategy under the equilibrium strategy for the restricted set of pure strategies at the current iteration  $(n_s)_j$ . We will then consider that the trained classifier is a close approximation to a best response to the equilibrium strategy of the steganographer. The major consequence of this heuristic is that *we lose the convergence guarantee of the double oracle algorithm*, as the algorithm is only guaranteed to converge if we compute an exact best response for both players at each stage.

## IV. $\epsilon$ -REGRET SOLUTIONS

Despite the relative efficiency of the double oracle algorithm compared to the standard method of linear programming for

**Algorithm 1** Double Oracle Algorithm**Data:** $n_c$ : Distribution of cover sources $\mathbf{X}$ : Pure strategies of the steganographer, $\mathbf{Y}_1$ :  $(f_1, \tau_1) \in \mathcal{R}_f$ , $\epsilon > 0$ **Result:**  $(f^*, \tau^*, n_s^*)$ :  $\epsilon$ -Nash equilibrium**repeat**Compute  $(f_j^*, \tau_j^*, (n_s^*)_j)$ , equilibrium of  $(\mathbf{X}, \mathbf{Y}_j, u)$ ;Compute  $(n_s)_{j+1} = \arg \max_{n_s \in \mathbf{X}} u(f_j^*, \tau_j^*, n_s)$ ;  
Train  $f_{j+1}$  using a training set with the proportion of cover and stego sources given by  $n_c$  and  $(n_s^*)_j$ ;Compute  $\tau_{j+1} = \arg \min_{\tau \in \bar{\mathbb{R}}} u(f_{j+1}, \tau, (n_s^*)_j)$  $\mathbf{Y}_{j+1} \leftarrow \mathbf{Y}_j \cap (f_{j+1}, \tau_{j+1})$ **until**  $u(f_j^*, \tau_j^*, (n_s)_{j+1}) - u(f_{j+1}, \tau_{j+1}, (n_s^*)_j) \leq \epsilon$ ;

solving large two-player zero-sum games, it is still very costly in practice since it requires training a new classifier at each iteration. Furthermore, since the number of iterations is not known a priori, the total number of classifiers that will need to be trained to obtain an acceptable solution is also unknown, and likely quite large. To solve this problem, we construct a class of so-called “ $\epsilon$ -regret” classifiers and show sufficient conditions under which they lead to an  $\epsilon$ -Nash equilibrium while also being guaranteed to necessitate only to train as many classifiers as cover sources plus one.

The basic idea of this section is to construct a classifier which is an  $\epsilon$ -best response to the best strategy of the steganographer. One way to do this is to find a dominant strategy, unavailable to the steganalyst in practice. This strategy will allow us to find a lower-bound to the game’s utility  $u$  that the steganalyst will try to approximate. Note that the steganalyst wants to construct a single classifier in order to ensure lower computational complexity than the double oracle algorithm. This means that only pure strategies are allowed for the steganalyst, contrary to the double oracle which leverages mixed strategies by randomizing many classifiers.

**A. Nash Equilibrium With Perfect Source Identification**

In our case, the most natural way to find a dominant strategy for the game defined by  $u$  is to assume that the steganalyst is able to perfectly identify each source. We will thus construct a lower-bound on the utility  $u$  by finding a Nash equilibrium in a game where the steganalyst is able to select a different classifier for each cover source.

To facilitate the discussion we assume the following:

- For each source  $i$ , there is a unique classifier  $f_i^*$  which, for all thresholds  $\tau \in \bar{\mathbb{R}}$ , is most powerful with respect to every other  $f \in \mathcal{K}$  when discriminating between  $P_i^c$  and  $P_i^s$ . Such a classifier can be thought of as a likelihood ratio test or more practically a classifier trained specifically on a target source.

- The ROC curve of each classifier  $f \in \mathcal{K}$  is concave and continuously differentiable on  $[0, 1]$ .
- The output distributions of  $f_i^*(x)$  has the same support when  $x \sim P_i^c$  and  $x \sim P_i^s$ .

Since the best classifier for each source is known, the goal of the steganalyst is simpler: they only have to find the threshold  $\tau_i$  for each source in order to minimize the utility of the steganographer. Indeed, since there is a single most powerful classifier for each source, and since the steganalyst can choose a different classifier for each source, the false alarms and missed detection errors are minimized by choosing  $f_i^*$ .<sup>1</sup>

To ease the presentation, we assume that, for a given classifier  $f_i^*$ , the steganalyst selects a probability of missed detection  $\beta_i^{f_i^*}$  and obtains the corresponding minimum probability of false alarm for this  $\beta_i^{f_i^*}$ :

$$\alpha_i^f(x) = \min_{\tau \in \bar{\mathbb{R}}} \{\alpha_i^f(\tau) \mid \beta_i^f(\tau) = x\}. \quad (17)$$

Note that we will drop the  $(x)$  from  $\alpha_i^f(x)$  when it is clear from context.

This game, which we will refer to as the *atomistic game*, is characterized by a utility  $u_A$  which we define as:

$$u_A(\beta, n_s) \triangleq \sum_{i=1}^N \lambda \alpha_i^{f_i^*} n_{c_i} + \bar{\lambda} \beta_i^{f_i^*} n_{s_i}. \quad (18)$$

Consequently, if there exists a pure strategy equilibrium for the steganalyst, it is a minimax strategy solving:

$$\begin{aligned} \max_{n_s \in \mathcal{R}_s} \min_{\beta \in [0,1]^N} u_A(\beta, n_s) \\ \text{s.t. } \sum_{i=1}^N n_s = 1. \end{aligned} \quad (19)$$

Notice that  $u(\beta, n_s)$  is convex in  $\beta$  (by concavity of the ROC curve) and concave (linear) in  $n_s$ . Hence, if a feasible solution exists, it obtained by finding the stationary point of the Lagrangian:

$$\mathcal{L} = u_A(\beta, n_s) + \eta \left( \sum_{i=1}^N n_{s_i} - 1 \right), \quad (20)$$

where  $\eta$  is the Lagrange multiplier. We have:

$$\frac{\partial \mathcal{L}}{\partial \beta_i} = 0 \Leftrightarrow n_{s_i} = - \frac{\partial \alpha_i}{\partial \beta_i} n_{c_i} \lambda \bar{\lambda}^{-1}, \quad (21)$$

$$\frac{\partial \mathcal{L}}{\partial n_{s_i}} = 0 \Leftrightarrow \beta_i^{f_i^*} = \eta \bar{\lambda}^{-1}. \quad (22)$$

Observe that at the equilibrium, Eq. (21) tells us that, the steganalyst should select a threshold such that every classifier has the same probability of missed detection. This is expected since such a choice guarantees that any convex sum of the probabilities of missed detection will always have the same value. Therefore, the steganographer is guaranteed to be unable

<sup>1</sup>Except in the case where an LRT is available, different classifiers might be optimal at different regimes. For example, at very low false-alarm rates, specific classifiers might have to be trained to obtain good results. Our abstraction still holds in this case as we can construct a function which associates to each threshold the “optimal classifier” for this threshold.

to increase her utility. Regarding the steganographer's strategy at the equilibrium, the probability that the steganographer embeds in the  $i$ -th source,  $n_{s_i}$ , not only depends on the environment through the probability of drawing a cover from the  $i$ -th source,  $n_{c_i}$ , but also on the classifiers available to the steganalyst through the slope of the ROC curve  $\frac{\partial \alpha_i}{\partial \beta_i}$ . Therefore the intuition that a steganographer should try to imitate the environment is not valid in the case where the steganographer has complete information over the steganalyst available strategies. Also, note that if we assume bounded  $\frac{\partial \alpha_i}{\partial \beta_i}$ , we have that if  $n_{c_i} = 0$  then  $n_{s_i} = 0$  at the equilibrium. This reflects the fact that if a cover source is absent from the environment, then the steganalyst can set their probability of false alarm to 1 on this cover source, thus being guaranteed to always detect any stego images coming from this source without incurring any cost in terms of false-alarm. On the other hand, under the hypothesis that the support of the cover and stego distribution is equal, it is impossible for the derivative  $\frac{\partial \alpha_i}{\partial \beta_i}$  to be equal to zero. Thus, as long as  $n_{c_i} > 0$ , the steganographer should always set  $n_{s_i}$  to be strictly positive, no matter how easy the corresponding cover source is to steganalyze.

The system in Eq (19) is easily solved by a bisection search on  $\beta$ . However, this system might not have a feasible solution in the first place, which would imply that there is no pure strategy for the steganalyst which leads to a Nash equilibrium.

The following proposition provides a sufficient condition for the existence (and uniqueness) of a solution to Eq (19):

*Proposition 1:* Let  $\lambda^* = \min(\lambda, \bar{\lambda})$ . Under the assumptions of this subsection, if we have that for all atomistic classifiers  $f_i^*$ :

$$\alpha_i^{f_i^*}(\lambda^*) \leq \lambda^*, \quad (23)$$

then the constrained optimization problem defined in Eq (19), admits a solution  $(\beta^*, n_s^*)$ . In particular it is always the case for  $\lambda = 0.5$ .

*Proof:* See Appendix A.  $\square$

In the rest of the paper, we assume an equilibrium solution exists and denote it as  $(\beta^*, n_s^*)$ . It is important to note that the atomistic utility  $u_A(\beta^*, n_s^*)$  is a lower-bound of any utility of the original game where the cover sources cannot be identified perfectly:

$$\forall (f, \tau) \in \mathcal{R}_f, u_A(\beta^*, n_s^*) \leq u(f, \tau, n_s^*). \quad (24)$$

*Atomistic game with two sources:* We end the discussion of the perfect source identification case by studying a simple example of a game containing only two cover sources. We assume that  $\lambda = 0.5$  and that the output of each optimal classifier  $f_i^*$  is Gaussian such that:

$$x \sim P_i^c \Rightarrow f_i^*(x) \sim \mathcal{N}(0, 1), \quad (25)$$

$$x \sim P_i^s \Rightarrow f_i^*(x) \sim \mathcal{N}(\mu_i, 1), \quad (26)$$

where  $\mu_i$  is the mean shift due to the embedding for source  $i$ . We want to note that such a setting closely matches the asymptotic distributions of the output of a likelihood ratio test in a simple hypothesis test such as the one found in MiPOD [19].

In such a case, a Nash equilibrium is found by solving the following system for  $\beta$ :

$$\begin{cases} n_{s_i} &= \exp\left(-Q^{-1}(\beta)\mu_i - \frac{\mu_i^2}{2}\right)n_{c_i} \\ 1 &= n_{s_1} + n_{s_2}. \end{cases} \quad (27)$$

where  $Q^{-1}$  is the inverse of the tail function of the standard Gaussian distribution. We show in Figure 2 how the solution of the game changes for different environments and source difficulties.

First of all, observe that in the case where the cover source environments are balanced –  $n_c = (0.5, 0.5)$  – the figure are symmetrical as expected. Indeed, this symmetry means that neither the steganographer nor the steganalyst favor one source over the other if they are identical. Furthermore, the behavior of each player is intuitive in this balanced environment. For the steganographer, the easier a source is relative to the other the less she will use it to create stego objects, i.e if  $\mu_1 > \mu_2$  then  $n_{s_1}^* < n_{s_2}^*$ . Similarly, for the steganalyst, each cover source contributes equally to a change in  $\beta^*$ .

In the case of an unbalanced environments, the behavior of the steganographer is less intuitive. Indeed, the steganographer will distribute most of the stego images in a cover source which is far easier to steganalyze than the other if this cover source is more prevalent. For example, look at the case where  $n_c = (0.8, 0.2)$ , that is where 80% of the environment consists of the first cover source. We observe that only when embedding in the first cover source becomes highly detectable – around  $\mu_1 = 4$  – and only when embedding in the second source is relatively safe – around  $\mu_2 = 1.5$ , then, and only then, will the steganographer prefer the second source over the first. Similarly, we can also observe that the less prevalent a source is, the less it impacts the choice of probability of missed detection  $\beta^*$  for the steganalyst – see that for  $n_c = (0.9, 0.1)$ , the contours are almost independent of  $\mu_2$ .

## B. Minimizing $P_E$ Over Atomistic Classifiers

We are in a position where we know how to construct a Nash equilibrium (up to the existence of a pure feasible solution) in the case where the steganalyst is able to perfectly identify the cover source of images. We can now go back to the original problem which does not assume access to such information. To solve this problem without incurring the high cost of the double oracle algorithm, we propose to train a single classifier which tries to match the lower-bound found in Section IV-A.

In order to design this classifier, we must find a suitable loss function to minimize. As a first step, we will review the loss function proposed in [16] and show why it is not suitable for our purpose. The idea of the classifier in [16], which we will denote  $f^\epsilon$ , is to minimize the maximum regret over all sources, the regret being computed with respect to the  $P_E$  of each  $f_i^*$ . Let:

$$\beta_i^{P_E} \triangleq \arg \min_{\beta \in [0,1]} \lambda \alpha_i^{f_i^*}(\beta) + \bar{\lambda} \beta, \quad (28)$$



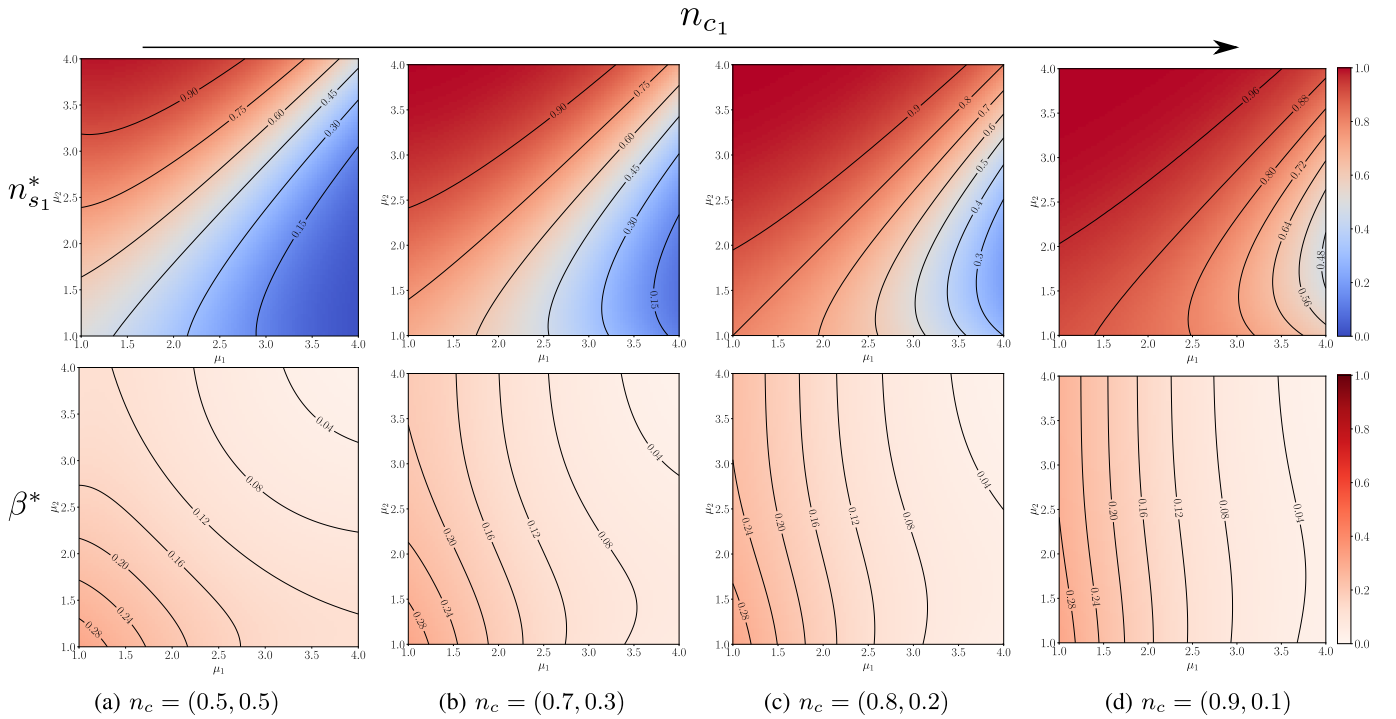


Fig. 2. Solutions  $(\beta^*, n_s^*)$  of the two-sources game as a function of the sources' "difficulty" given by the distribution shift  $(\mu_1, \mu_2)$ . The higher  $\mu_i$ , the easier the source is to steganalyze. The stego-source distribution of the first source  $n_{s_1}^*$  is displayed at the top and the probability of missed detection of the classifier at the equilibrium  $\beta^*$  at the bottom. For the top figures, a blue color means the steganographer favors the second cover source whereas a red color means she favors the first one.

be the probability of missed detection corresponding to the  $P_E$  for each atomistic classifier.<sup>2</sup>

The regret of [16] can then be defined as:

$$\epsilon \triangleq \min_{f \in \mathcal{K}} \max_{1 \leq i \leq N} \text{Regret}(f, i), \quad (29)$$

$$\text{Regret}(f, i) = \lambda \left( \alpha_i^f - \alpha_i^{f_i^*}(\beta_i^{P_E}) \right) + \bar{\lambda} \left( \beta_i^f - \beta_i^{P_E} \right).$$

and thus  $f^\epsilon \triangleq \arg \min_{f \in \mathcal{K}} \max_{1 \leq i \leq N} \text{Regret}(f, i)$ .

In other words, the min-regret classifier proposed in [16] tries to match the  $P_E$  of each atomistic classifier  $f_i^*$  with a single classifier and a single threshold. Note that, in order to ease the notation, we assume the threshold of each  $f$  to always be set to the optimal one and thus omit it in the notation.

Using the knowledge of Section IV-A, we can first question the choice of using the  $P_E$  as a score to match. Indeed, if we reach a regret of 0 under the definition of Eq (29), we would not be guaranteed to reach a strategy that leads to a Nash equilibrium in the atomistic case. Indeed, Eq (22) tells us that we need each  $\beta_i^{P_E}$  to be equal to the same value, which is highly unlikely except in the trivial case where all sources are identical. Consequently, in general, there does not exist a  $\beta^*$  such that  $\beta_i^{P_E} = \beta^*$  for all sources, which would mean that  $\beta_i^{P_E}$  is not an equilibrium strategy for the steganalyst in the atomistic case, hence not a dominant strategy in our current setting.

Nevertheless, it is simple to fix this design choice by matching the atomistic equilibrium error  $\lambda \alpha_i^{f_i^*}(\beta_i^*) + \bar{\lambda} \beta_i^*$  for

each  $f_i^*$  instead of the  $P_E$ . This leads us to the following re-definition of the regret:

$$\epsilon \triangleq \min_{f \in \mathcal{K}} \max_{1 \leq i \leq N} \text{Regret}(f, i), \quad (30)$$

$$\text{Regret}(f, i) = \lambda \left( \alpha_i^f - \alpha_i^{f_i^*}(\beta_i^*) \right) + \bar{\lambda} \left( \beta_i^f - \beta_i^* \right).$$

We will now show that even by matching the atomistic equilibrium error of each atomistic classifier  $f_i^*$ , we are still not guaranteed to reach an  $\epsilon$ -Nash equilibrium.

Let us first define the differences of false-alarms and missed detections between the min-regret classifier and the atomistic classifiers as:

$$\epsilon_{\alpha_i} = \lambda \left( \alpha_i^{f^\epsilon} - \alpha_i^{f_i^*}(\beta_i^*) \right), \quad (31)$$

$$\epsilon_{\beta_i} = \bar{\lambda} \left( \beta_i^{f^\epsilon} - \beta_i^* \right). \quad (32)$$

We can now compare the performance of the min-regret classifier  $f^\epsilon$  to  $f^*$ :

$$u(f^\epsilon, n_s^*) - u_A(\tau^*, n_s^*) = \sum_{i=1}^N \epsilon_{\alpha_i} n_{c_i} + \sum_{i=1}^N \epsilon_{\beta_i} n_{s_i}^*. \quad (33)$$

It is straightforward to construct examples where the min-regret classifier is not guaranteed to lead to an  $\epsilon$ -best response to  $n_s^*$ . For example, assuming that for all  $i$ ,  $n_{s_i}^* \neq 0$

<sup>2</sup>Note that to ease the presentation, we still call this quantity  $P_E$  even though we might have  $\lambda \neq 0.5$ .

and  $n_{c_i} \neq 0$ , let us fix  $\epsilon_{\alpha_i}$  as:

$$\epsilon_{\alpha_i} = \begin{cases} -\frac{\epsilon}{n_{c_i} n_{s_i}^*} & \text{if } n_{s_i}^* \geq n_{c_i} \\ \frac{\epsilon}{n_{c_i} n_{s_i}^*} & \text{if } n_{s_i}^* < n_{c_i} \end{cases} \quad (34)$$

and similarly for  $\epsilon_{\beta_i}$ :

$$\epsilon_{\beta_i} = \begin{cases} \epsilon + \frac{\epsilon}{n_{c_i} n_{s_i}^*} & \text{if } n_{s_i}^* \geq n_{c_i} \\ \epsilon - \frac{\epsilon}{n_{c_i} n_{s_i}^*} & \text{if } n_{s_i}^* < n_{c_i} \end{cases} \quad (35)$$

Then we indeed have for all  $i$  that  $\epsilon_{\alpha_i} + \epsilon_{\beta_i} = \epsilon$ , yet we also have:

$$\epsilon_{\alpha_i} n_{c_i} + \epsilon_{\beta_i} n_{s_i}^* = \begin{cases} \epsilon n_{s_i}^* + \frac{\epsilon}{n_{c_i}} - \frac{\epsilon}{n_{s_i}^*} & \text{if } n_{s_i}^* \geq n_{c_i} \\ \epsilon n_{s_i}^* + \frac{\epsilon}{n_{s_i}^*} - \frac{\epsilon}{n_{c_i}} & \text{if } n_{s_i}^* < n_{c_i} \end{cases} \quad (36)$$

$$\triangleq \epsilon n_{s_i}^* + \epsilon'_i,$$

where  $\epsilon'_i \geq 0$ , with equality only when  $n_{c_i} = n_{s_i}^*$ .

This leads us to the actual utility difference:

$$\sum_{i=1}^N \epsilon_{\alpha_i} n_{c_i} + \epsilon_{\beta_i} n_{s_i}^* = \epsilon + \sum_{i=1}^N \epsilon'_i \geq \epsilon. \quad (37)$$

Therefore, the min-regret classifier  $f^\epsilon$  is not, in general, guaranteed to be an  $\epsilon$ -best response to the equilibrium strategy of the steganographer  $n_s^*$ . Note here that even if we replaced atomistic classifier  $f_i^*$  by the actual best response to  $n_s^*$  available to the steganalyst in  $\mathcal{R}_f$ , we could construct a counter-example using the same method to show that this definition of regret does not lead to an  $\epsilon$ -best response.

Now we also need to know when  $n_s^*$  is an  $\epsilon$ -best-response to  $f^\epsilon$ . We know that if the steganalyst plays the min-regret classifier  $f^\epsilon$ , then the best response of the steganographer is to set  $n_{s_i} = 1$  where  $i = \arg \max_i \beta_i^{f^\epsilon}$ . Hence for  $n_s^*$  to be an  $\epsilon$ -best-response to  $f^\epsilon$  we need:

$$\left( \sum_{i=1}^N \bar{\lambda} \beta_i^{f^\epsilon} n_{s_i}^* \right) - \max_i \bar{\lambda} \beta_i^{f^\epsilon} \geq -\epsilon, \quad (38)$$

which is, once again, not true in general.

In summary, the naive definition of regret as minimizing the bound on the  $P_E$  individually for each atomistic classifier fails to provide a  $\epsilon$ -Nash equilibrium. More generally, even minimizing such a bound on the ‘‘correct’’ sum of errors does not provide any guarantees. The next subsection’s goal is to provide another definition that guarantees such an equilibrium while still being simple to compute in practice.

### C. Min $\epsilon$ -Regret Classifier

As we have seen, minimizing the maximum sum of errors individually over sources is too weak to actually guarantee an  $\epsilon$ -Nash equilibrium. In particular, we have seen that we need to take into account:

- The impact of the environment on the sum of  $\epsilon_\alpha$  and  $\epsilon_\beta$  – c.f Eq (33),

- The maximum  $\beta_i^{f^\epsilon}$  in order to control the steganographer’s best response – c.f Eq (38)

For all classifiers  $f \in \mathcal{K}$ , we can define the regret of the steganalyst using Eq (33) and the regret of the steganographer using Eq (38):

$$R_{\text{def}}(f) = \sum_{i=1}^N \epsilon_{\alpha_i} n_{c_i} + \epsilon_{\beta_i} n_{s_i}^*, \quad (39)$$

$$R_{\text{atk}}(f) = \max_i \bar{\lambda} \beta_i^f - \left( \sum_{i=1}^N \bar{\lambda} \beta_i^f n_{s_i}^* \right). \quad (40)$$

This leads us to a natural redefinition of our min-regret classifier  $f^\epsilon$  in the case of our game:

$$\epsilon \triangleq \min_{f \in \mathcal{K}} \text{Regret}(f), \quad (41)$$

$$\text{Regret}(f) = \max(R_{\text{def}}, R_{\text{atk}}). \quad (42)$$

Following Eq (33) and Eq (38), the min-regret classifier  $f^\epsilon$  is guaranteed to be an  $\epsilon$ -Nash equilibrium.

In order to avoid the complexity of the double-oracle algorithm, we can thus train a single classifier using regret as defined in (41) as a loss function. In particular, such a loss can be used by any deep neural network. However, it is not directly usable since the computation of both the probabilities of false-alarm  $\alpha_i^f$  and of the probabilities of missed detection  $\beta_i^f$  are non-differentiable due to the thresholding operation.<sup>3</sup>

This is a well-known problem in the literature and is usually solved by replacing the step function with a surrogate function which is differentiable and sometimes convex [20]. In our case, we will simply use a sigmoid function on the soft outputs of our neural net:

$$\sigma(x) = \frac{e^{r(x-0.5)}}{1 + e^{r(x-0.5)}}. \quad (43)$$

In order to make the optimization efficient we also need good estimates of the probabilities of missed detection  $\beta_i^f$  at each step in order to select the correct maximum. This can be accomplished by using a large batch size but this strategy would not scale with the number of sources. In order to make our method scalable, we retain the soft outputs of each mini-batch and use them to compute the  $\beta_i^f$  at each optimization step.

## V. EXPERIMENTS

We now end the paper by studying the performance of the three proposed methods: the atomistic strategy, the double oracle strategy and the min-regret strategy. In particular, we provide a comparison against the standard holistic classifier which does not assume a rational steganographer in order to show the importance of this assumption.

### A. Experimental Core Settings

All experiments were performed using a dataset of 1150 RAW images from the ALASKA dataset [6]. These images were all taken by a Canon EOS 100D camera with

<sup>3</sup>The max operation is itself not differentiable but modern machine learning packages such as *Pytorch* solve this problem automatically by using the sub-gradient instead.

	Source 1	Source 2	Source 3	Source 4	Source 5	Source 6	Source 7	Source 8	Source 9
Source 1	1.6	+0.6	+1.4	+6.1	+9.9	+7.6	+7.4	+6.0	+2.2
Source 2	+1.1	4.7	+1.9	+2.3	+9.1	+6.7	+7.2	+4.8	+3.2
Source 3	+1.3	+0.8	7.3	+2.4	+12.6	+6.9	+4.4	+3.0	+3.0
Source 4	+8.4	+4.6	+3.2	13.1	+16.3	+6.0	+5.3	+5.7	+3.4
Source 5	+1.2	+0.5	+3.1	+4.4	18.9	+6.2	+5.9	+1.4	+1.9
Source 6	+15.1	+13.4	+10.1	+13.1	+27.9	29.3	+3.2	+11.8	+5.3
Source 7	+15.3	+18.1	+12.8	+18.5	+26.2	+6.6	34.3	+11.7	+5.6
Source 8	+5.1	+3.2	+2.6	+3.7	+5.2	+9.0	+6.3	35.4	-0.2
Source 9	+15.1	+13.5	+8.3	+9.2	+13.8	+11.4	+7.8	+2.8	43.1
Holistic	+0.4	-0.7	-0.2	-0.2	+3.0	+3.0	+0.5	+0.9	+0.7

Fig. 3.  $P_E$  and source inconsistency for the 9 sources used in the paper. Steganography performed using JMIPOD at 0.2 bpc and steganalysis using EfficientNet-b3. Rows correspond to the source used in the training set while columns correspond to the testing set. The “holistic” row corresponds to the classifier trained on a dataset composed of all 9 sources. The diagonal, in green, corresponds to the  $P_E$  where there is no cover-source mismatch (intrinsic difficulty). The off-diagonal elements correspond to the difference of the  $P_E$  of the mismatched classifier with the  $P_E$  of the matched classifier in the same column (source inconsistency).

ISO ranging from 1000 to 6400. From these RAW images, a cropping operation was performed to obtain 13,380 non-overlapping images of size  $264 \times 264$ . Finally, the dataset was processed with different processing pipelines using the method described in [17, Fig. 1] leading to 243 cover source image sets, of which 9 were selected as representatives using the method given in Section III from the same work. This led to 9 cover source image sets of  $264 \times 264$  JPEG images at quality factor 98. Each of them was split as 10380/1000/2000 for training/validation/testing.

All steganography was performed using JMIPOD [21] which is the current state-of-the-art of non side-informed JPEG steganography.

Steganalysis was always performed using a pretrained EfficientNet-b3 [22] which is the current state-of-the-art in steganalysis as shown during the ALASKA2 competition [1] as well as in the study of Yousfi et al. [2], [3]. Atomistic classifiers, which will be considered as our  $f_i^*$ , were trained using curriculum learning starting at 0.4 bpc for 30 epochs directly down to 0.2 bpc for 15 epochs. The learning rate was set to 0.5 and divided by two on loss plateau and the batch-size set to 64. Holistic detectors were also trained by using a training set containing 10380 images with an equal number of image per source. The best models for all classifiers was selected by taking the model with the lowest loss on the validation set.

We present the performance of the atomistic and holistic classifiers in Figure 3. As expected, *sources have highly different intrinsic difficulties* – the  $P_E$  in the ideal case when the source of the test set matches the source of the atomistic classifier – going from 1.6% up to 43%. This demonstrates the diversity of the cover source environment as would be expected in a real-world environment. Furthermore, the *inconsistency* between sources – the difference of  $P_E$  when the source of the testing set does not match the source of the atomistic classifier – goes from 0.5% to 27% with an average of 6.5% for each column. On the other hand, notice that the  $P_E$  of the holistic classifier is almost always very close to the  $P_E$  of the atomistic classifiers with an average of only 0.8% of difference. This begs the question: is the holistic classifier

TABLE I  
AVERAGE GAME VALUES AND CLOSENESS TO A NASH EQUILIBRIUM FOR JMIPOD AT 0.2 BPC AND EFFICIENTNET-B3

	$u_{atk}(f, \tau, n_s^*)$	$u(f, \tau, n_s^*)$	Difference from Nash equilibrium
$f^*, \beta^*$	17.6%	24.5%	-
$f^{hol}$	37.6%	32.8%	-
$f^{DO}$	16.5%	26.2%	0.093%
$f^\epsilon$	11.2%	25.4%	$\leq 2.2\%$

itself sufficient to obtain a good approximation of the best response of the steganalyst when perfect source identification is available? We will answer this question by comparing this holistic classifier to the different method we have presented in the paper.

### B. Comparison of Proposed Solutions

In order to evaluate the performance of our proposed solutions we compare the utility of the game when using:

- The atomistic optimal strategy as defined in Section IV-A, denoted as  $(f^*, \beta^*)$ ,
- The holistic classifier, denoted  $f^{hol}$ , serving as the state-of-the-art baseline since it is the current strategy of choice for heterogeneous environments [1], [5], [6],
- The double oracle solution as defined in Section III-B, denoted  $f^{DO}$ ,
- The minimum regret solution as defined in Section IV-C, denoted  $f^\epsilon$ .

When using the double oracle algorithm, we compute the steganalyst best-response by first initializing the weights of EfficientNet to those of the holistic classifier in order to speed up convergence. The training is then run only for 5 epochs for each step of the double oracle algorithm. The double oracle algorithm was always run for 25 steps.

We perform the same initialization when computing the min-regret solution and run the training for 15 epochs as no gain was observed beyond this point. The best model is selected as the model with lowest regret on the validation set. The sigmoid parameter was set to  $r = 25$ . This value was selected by a simple grid search between 1 and 50 minimizing the regret on a random environment of cover sources. It should be noted that the exact value of  $r$  did not significantly impact the convergence during the training step – see Figure 5.

Finally, all experiments were performed assuming  $\lambda = 0.5$  in order to guarantee the existence of a solution for the case when sources can be perfectly identified. We would like to emphasize that this means that the probability of false alarm and of missed detection are equally weighted, *which makes the zeros-sum game utility,  $u$ , close to the usual  $P_E$  metric used to evaluate steganography and steganalysis.*

The evaluation was performed by sampling randomly 10 different cover source environments distributions  $n_c$ . We provide the average probability of error  $u(f, \tau, n_s^*)$  for each classifier and how close they are to Nash equilibrium in Table I. Note that in the case of the original holistic classifier  $f^{hol}$ , we assume that the steganographer plays the best response against it.

	Source 1	Source 2	Source 3	Source 4	Source 5	Source 6	Source 7	Source 8	Source 9	All sources
$P_E$	1.6	4.7	7.3	13.1	18.9	29.3	34.3	35.4	43.1	22.1
$n_c$	21.3	4.4	6.0	24.7	2.2	6.4	7.7	8.5	18.8	-
$\alpha^{f^*}, \tau^*$	0.0	0.2	1.6	15.3	23.8	38.7	47.6	59.4	75.1	29.8
$\alpha^{f^{hol}}$	+1.6	+2.3	+4.4	-4.2	-2.3	-2.3	-14.4	-11.0	-7.6	-4.0
$\alpha^{f^\epsilon}$	+1.5	+3.4	+8.1	+1.6	+10.8	+22.4	+15.9	+5.8	+11.1	+6.8
$\alpha^{f^{DO}}$	+0.3	+4.7	+5.8	+4.1	+20.6	+19.8	+14.8	+7.1	+2.1	+5.5
$n_s^*$	0.0	0.1	1.3	25.1	2.6	7.9	11.1	13.7	38.1	-
$\beta^{f^*}, \tau^*$	16.6	16.6	16.6	16.6	16.6	16.6	16.6	16.6	16.6	16.6
$\beta^{f^{hol}}$	-13.7	-9.9	-7.6	-0.7	+6.2	+11.1	+21.0	+9.6	+4.1	+21.0
$\beta^{f^\epsilon}$	-14.0	-11.7	-10.4	-5.6	-3.4	-8.4	-5.0	-5.5	-11.7	-3.4
$n_s^{DO}$	0.0	0.0	0.0	31.1	6.9	18.5	12.0	9.2	22.4	-
$\beta^{f^{DO}}$	-5.3	-8.1	-5.6	-5.1	-5.1	-5.1	-5.1	-5.1	-5.1	-5.1

Fig. 4. Probability of false alarm ( $\alpha$ ) and probability of missed detection ( $\beta$ ) for each source depending on the classifier used. The  $P_E$  is computed using the atomistic classifier  $f_i^*$  for individual sources. When computing the  $P_E$  for the dataset containing all sources, we assume a balanced dataset in terms of sources.  $f^{hol}$  corresponds to the holistic classifier trained on all sources,  $f^\epsilon$  to the min  $\epsilon$ -regret solution and  $f^{DO}$  to the double oracle solution. Note that  $n_s^{DO}$  corresponds to the  $\epsilon$ -best response of the steganographer when  $f^{DO}$  is played. Note that for computing  $\beta^{f^{hol}}$  and  $\beta^{f^\epsilon}$  for the dataset containing all sources, we assume the steganographer plays a best response against these classifiers and not  $n_s^*$  – optimal only for  $f^*$ . The sources are sorted by  $P_E$  in ascending order. Blue cells correspond to environments, white cells to the optimal performance on the atomistic game. Cells are red (resp. green) when the performance of the classifier is worse (resp. better) than the performance of the solution for the atomistic game.

As a first observation, note that simply using the holistic classifier is not a good response for the steganalyst: the absolute difference with respect to the atomistic game value is on average equal to 8.7%. This means that either the steganographer can choose a cover source to significantly degrade the steganalyst's performance. *It demonstrates that a classifier with  $P_E$  close to the atomistic classifiers for every source is not guaranteed to be close to a best response for the steganalyst.*

This leads us to the proposed methods, both the min-regret method and the double oracle algorithm provide solutions with values close to the game where perfect identification of sources is possible while also being close to a Nash equilibrium. From the results, it is clear that the double oracle algorithm provides far better guarantees than the min-regret method with a solution guaranteed to be within 0.1% of a Nash equilibrium whereas our method only guarantees the solution to be **at worst** 2.2% from a Nash equilibrium. However it should be understood that the closeness to a Nash equilibrium, denoted as  $\epsilon$ , is computed differently for the double oracle and the min-regret. In the case of the double oracle,  $\epsilon$  is computed as  $u(f_j^*, \tau_j^*, (n_s)_{j+1}) - u(f_{j+1}, \tau_{j+1}, (n_s^*)_j)$  at each step, hence providing a *guarantee without reference to the atomistic game*  $u_A$ . On the other hand, the  $\epsilon$  of the min-regret solution is computed by comparing it to the atomistic game solution as  $u(f^\epsilon, \tau, n_s^*) - u_A(\beta^*, n_s^*)$ . However, this solution only provides a lower-bound for the steganalyst – and hence an upper-bound on the lost utility – which might not be achievable since perfect identification of the sources is not available.

It might then be possible that the min-regret solution is actually a lot closer to the available pure best-response than what the regret is actually guaranteeing. Once again, it should also be noted that despite slightly better solutions, the double oracle algorithm is far more costly with the need to compute a new classifier at each iteration, without any knowledge of how many iterations will allow to reach a suitable solution. On the other hand, the min-regret method needs only as many classifiers as sources – to compute the equilibrium for the atomistic game – plus one – the min-regret classifier itself. Nevertheless, both methods allow us to outperform the naive holistic classifier by at least 5.2% in the worst case, showing the importance of taking the cover source environment into account when the steganographer is assumed to be a rational adversary.

For a more precise evaluation of the different solutions, we present the detailed performance of each classifier for a single environment in Figure 4. First of all, observe that for both the double oracle and the atomistic case, the strategy of the steganographer follows the same trend. The leftmost sources, which are the easiest to steganalyse are almost not used, despite Source 1 making up 21% of cover environment. The weight of these “easy sources” is then transferred to more difficult sources. In particular, in the atomistic case, the steganalyst strategy  $n_s^*$  transfers these weights to the most difficult sources whereas the sources with intermediate difficulty keep the same proportion as for the covers. Another observation is that both the double oracle and min-regret solution trade-off higher probabilities of false-alarm for lower probabilities

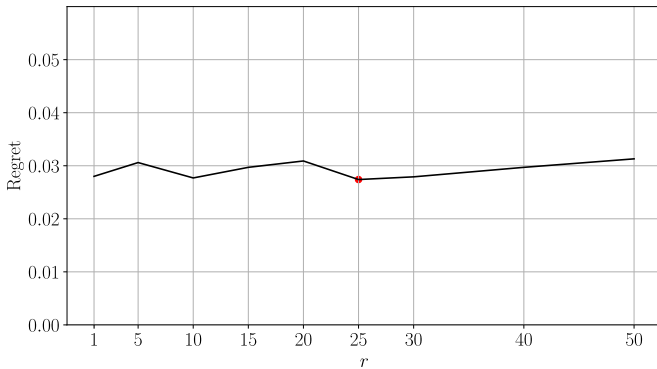


Fig. 5. Regret as a function of the softmax parameter  $r$  for JMIPOD at 0.2bpp on a random environment composed of 9 sources and steganalyzed with the min-regret classifier using EfficientNet-b3. The red marker indicates the minimum regret obtained at  $r = 25$ .

of missed detection. Finally notice that the atomistic strategy leads to a constant probability of false missed detection, as constrained by Eq (22). The double oracle solution also has an almost constant probability of missed detection across sources for the same reason.

Finally, regarding the training time of both algorithms, the computational complexity is in both cases dominated by the training of deep neural network. Under our hardware setting, using a single NVIDIA Tesla V100 32GB, a single epoch takes on average 5 min in both cases, leading to an average training time of 625 min for the double-oracle algorithm and only 75 min for the min-regret classifier. However, these results hold because of the fixed number of training steps we have set for both methods. If instead, we measure the training time such that the double-oracle reaches the best regret for the min-regret method, we measure a training time of 320 min for the double-oracle detector and 34 min for the min-regret detector. For both types of measurement, we thus observe that the min-regret method is approximately 10 times faster than the double oracle algorithm.

## VI. CONCLUSION

In this paper, we have proposed a formulation of the game played between the steganographer and the steganalyst as a two-player non-zero-sum game. This formulation naturally takes into account the presence of different cover sources and how it impacts the strategies of both the steganalyst and steganographer. We showed how to convert this non-zero-sum game into an equivalent zero-sum game. This allowed us to propose methods to compute solutions leading to a Nash equilibrium.

We showed that under the assumption that the steganalyst and the steganographer are both rational, the steganalyst should select a threshold such that the classifier has identical probability of missed detection on every cover source. On the other hand, the steganographer should still use “easy to steganalyze” cover sources if they are prevalent. In particular, she should never forego using a cover source – i.e  $n_{c_i} > 0 \Rightarrow n_{s_i} > 0$ . Furthermore, we have shown that both the double oracle algorithm and our min-regret classifier allowed reaching solutions close to a Nash equilibrium. The advantage of the

min-regret classifier being its extremely low computational complexity compared to the double oracle algorithm, at the cost of slightly weaker guarantees in terms of closeness to a Nash equilibrium. Finally, the standard holistic classifier was shown to be extremely vulnerable to a rational steganographer: chasing only the performance of atomistic classifiers in a realistic context should thus be considered a mistake for the steganalyst. We emphasize that our min-regret method can be applied to **any** detector based on loss function in order to make it robust to a rational steganographer.

The model we used in this work could be extended in several directions. In particular, we assumed that the steganographer uses only a single stego algorithm with a fixed payload. Furthermore, we have assumed both players to be fully rational and to have complete knowledge of the environment. However, it can be argued that the steganalyst has more resources than the steganographer and thus should be able to obtain better estimates of the environment. Future works should integrate these variables into the model to obtain a truly complete representation of the steganography game.

## APPENDIX

### A. Proof of Proposition 1

We will only treat the case where  $\lambda \leq 0.5$ , that is, when  $\lambda^* = \lambda$ , the other case is symmetrical.

Let  $g$  be a function defined as:

$$g : [0, 1] \rightarrow \mathbb{R}$$

$$\beta \mapsto \sum_{i=1}^N \alpha_i(\beta) n_{c_i}. \quad (44)$$

We have that:

$$g(0) = \sum_{i=1}^N n_{c_i} = 1, \quad g(1) = 0, \quad g(\lambda) \leq \lambda, \quad (45)$$

with the last inequality given by Eq (23).

By concavity of the ROC curves,  $g$  is convex, hence:

$$\frac{dg}{d\beta}(1) \geq -1 \geq -\bar{\lambda}\lambda^{-1}, \quad (46)$$

$$\frac{dg}{d\beta}(0) \leq \frac{g(\lambda) - g(0)}{\lambda} \leq \frac{\lambda - 1}{\lambda} = -\bar{\lambda}\lambda^{-1}. \quad (47)$$

Using the fact that the ROC curves are continuously differentiable on  $[0, 1]$ ,  $\frac{dg}{d\beta}$  is continuous on  $[0, 1]$ , hence by the intermediate value theorem we have that:

$$\exists \beta^* \in (0, 1), \quad \frac{dg}{d\beta}(\beta^*) = \sum_{i=1}^N \frac{\partial \alpha_i}{\partial \beta_i}(\beta^*) n_{c_i} = -\bar{\lambda}\lambda^{-1}. \quad (48)$$

Now, if we set  $(\beta^*, n_s^*)$  such that, for all  $i$ :

$$\beta_i^* = \beta^*,$$

$$n_{s_i}^* = -\frac{\partial \alpha_i}{\partial \beta_i}(\beta_i^*) n_{c_i} \bar{\lambda}^{-1}. \quad (49)$$

Then, we obtain a sufficient condition for the existence of a solution:

$$\sum_{i=1}^N n_{s_i}^* = 1 \iff \frac{dg}{d\beta}(\beta^*) = \bar{\lambda}\lambda^{-1} \iff \lambda \leq 0.5. \quad (50)$$

## REFERENCES

- [1] R. Cogranne, Q. Giboulot, and P. Bas, "ALASKA#2: Challenging academic research on steganalysis with realistic images," in *Proc. IEEE Int. Workshop Inf. Forensics Secur. (WIFS)*, New York, NY, USA, Dec. 2020, pp. 1–5.
- [2] Y. Yousfi, J. Butora, E. Khvedchenya, and J. Fridrich, "ImageNet pre-trained CNNs for JPEG steganalysis," in *Proc. IEEE Int. Workshop Inf. Forensics Secur. (WIFS)*, New York, NY, USA, Dec. 2020, pp. 1–6.
- [3] Y. Yousfi, J. Butora, J. Fridrich, and C. F. Tsang, "Improving EfficientNet for JPEG steganalysis," in *Proc. ACM Workshop Inf. Hiding Multimedia Secur.*, Jun. 2021, pp. 149–157.
- [4] Q. Giboulot, R. Cogranne, D. Borghys, and P. Bas, "Effects and solutions of cover-source mismatch in image steganalysis," *Signal Process., Image Commun.*, vol. 86, Aug. 2020, Art. no. 115888.
- [5] G. Quentin, B. Patrick, C. Rémi, and B. Dirk, "The cover source mismatch problem in deep-learning steganalysis," in *Proc. 30th Eur. Signal Process. Conf. (EUSIPCO)*, Aug. 2022, pp. 1032–1036.
- [6] R. Cogranne, Q. Giboulot, and P. Bas, "The Alaska steganalysis challenge: A first step towards steganalysis," in *Proc. ACM Workshop Inf. Hiding Multimedia Secur.*, Paris, France, Jul. 2019, pp. 125–137.
- [7] T. Pevný and A. D. Ker, "Towards dependable steganalysis," *Proc. SPIE*, vol. 9409, Mar. 2015, Art. no. 94090I.
- [8] A. D. Ker, T. Pevný, and P. Bas, "Rethinking optimal embedding," in *Proc. 4th ACM Workshop Inf. Hiding Multimedia Secur.* Vigo Galicia Spain, Jun. 2016, pp. 93–102.
- [9] P. Schöttle and R. Böhme, "Game theory and adaptive steganography," *IEEE Trans. Inf. Forensics Security*, vol. 11, no. 4, pp. 760–773, Apr. 2016.
- [10] S. Bernard, P. Bas, J. Klein, and T. Pevny, "Explicit optimization of min max steganographic game," *IEEE Trans. Inf. Forensics Security*, vol. 16, pp. 812–823, 2021.
- [11] H. B. McMahan, G. J. Gordon, and A. Blum, "Planning in the presence of cost functions controlled by an adversary," in *Proc. 20th Int. Conf. Mach. Learn.*, 2003, pp. 536–543.
- [12] T. Denemark, V. Sedighi, V. Holub, R. Cogranne, and J. Fridrich, "Selection-channel-aware rich model for steganalysis of digital images," in *Proc. IEEE Int. Workshop Inf. Forensics Secur. (WIFS)*, Dec. 2014, pp. 48–53.
- [13] V. Lisý, R. Kessl, and T. Pevný, "Randomized operating point selection in adversarial classification," in *Machine Learning and Knowledge Discovery in Databases*. Berlin, Germany: Springer, 2014, vol. 8725, pp. 240–255.
- [14] M. Barni and B. Tondi, "Theoretical foundations of adversarial binary detection," *Found. Trends Commun. Inf. Theory*, vol. 18, no. 1, pp. 1–172, 2020.
- [15] M. Barni, G. Cancelli, and A. Esposito, "Forensics aided steganalysis of heterogeneous images," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process.*, Dallas, TX, USA, Mar. 2010, pp. 1690–1693.
- [16] D. Šepák, L. Adam, and T. Pevný, "Formalizing cover-source mismatch as a robust optimization," in *Proc. EUSIPCO*, Belgrade, Serbia, 2022, pp. 1042–1046.
- [17] R. Abecidan, V. Itier, J. Boulanger, P. Bas, and T. Pevny, "Using set covering to generate databases for holistic steganalysis," in *Proc. IEEE Int. Workshop Inf. Forensics Secur. (WIFS)*, Shanghai, China, Dec. 2022, p. 6.
- [18] L. Adam, R. Horcik, T. Kasl, and T. Kroupa, "Double Oracle algorithm for computing equilibria in continuous games," in *Proc. AAAI Conf. Artif. Intelligence*, May 2021, vol. 35, no. 6, pp. 5070–5077.
- [19] V. Sedighi, R. Cogranne, and J. Fridrich, "Content-adaptive steganography by minimizing statistical detectability," *IEEE Trans. Inf. Forensics Security*, vol. 11, no. 2, pp. 221–234, Feb. 2016.
- [20] L. Adam, V. Mácha, V. Smídl, and T. Pevný, "General framework for binary classification on top samples," *Optim. Methods Softw.*, vol. 37, no. 5, pp. 1–32, Dec. 2021.
- [21] R. Cogranne, Q. Giboulot, and P. Bas, "Efficient steganography in JPEG images by minimizing performance of optimal detector," *IEEE Trans. Inf. Forensics Security*, vol. 17, pp. 1328–1343, 2022.
- [22] M. Tan and Q. V. Le, "EfficientNet: Rethinking model scaling for convolutional neural networks," 2019, *arXiv:1905.11946*.