

# Explaining the Bag Gain in Batch Steganography

Eli Dworetzky and Jessica Fridrich, *Fellow, IEEE*

**Abstract**—In batch steganography, the sender distributes the secret payload among multiple images from a “bag” to decrease the chance of being caught. Recent work on this topic described an experimentally discovered phenomenon, which we call the “bag gain”: for fixed communication rate, pooled detectors experience a decrease in statistical detectability for initially increasing bag sizes, providing an opportunity for the sender to gain in security. The bag gain phenomenon is universal in the sense of manifesting under a wide spectrum of conditions. In this paper, we explain this experimental observation by adopting a statistical model of detector response. Despite the simplicity of the model, it does capture observed trends in detectability as a function of the bag size, the rate, and cover source properties. Additionally, and surprisingly, the model predicts that in certain cover sources the sender should avoid bag sizes that are too small as this can lead to a bag loss.

**Index Terms**—Batch Steganography, Pooled Steganalysis

## I. INTRODUCTION

The problem of batch steganography and pooled steganalysis has been introduced by Ker in 2006 [16] and has since been a subject of intense research [17], [19], [22], [13], [20], [24], [26], [27], [25], [32], [31], [30]. Batch steganography deals with the situation when the sender spreads her payload among multiple covers (a bag of cover images) to decrease the Warden’s chances of detecting the use of this stealth communication channel. As formulated in the original work of Ker [16], if the steganographer is allowed to spread payload among multiple images, the Warden (or steganalyst) is free to pool evidence from the same multitude of images to detect the use of steganography, a process known as pooled steganalysis. In particular, the Warden uses a so-called pooled detector (or “pooler”) to decide whether a bag of images in its entirety is stego.

Intuitively, to improve security images that are harder to steganalyze should receive a larger payload and vice versa. Such batch senders have originally been studied in [26]. Most notably, the authors studied the so-called Image Merging Sender that assigns payloads to images from the bag by considering their union as a single larger image in which a message is embedded with some content-adaptive steganographic algorithm. This sender has also been studied in [27], [1] and most recently in [30]. In particular, our previous work [30] considered the scenario where both the sender and the Warden use feedback from a single-image steganography detector. Two

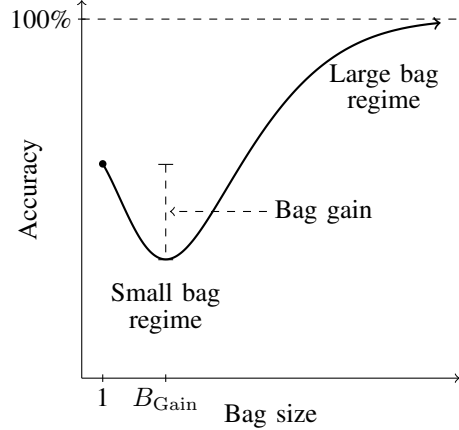


Figure 1. The universal trend of a pooled detector’s accuracy as a function of bag size  $B$  when a fixed positive communication rate is maintained. In the small bag regime, it is possible for the steganographer to gain security by spreading payload among  $B_{Gain}$  cover images. In the large bag regime, the detectability monotonically increases with bag size.

batch methods were introduced to exploit the information gain from a detector’s feedback and were compared with the detector-agnostic Image Merging Sender, bringing insight into the practical design of detector-informed senders and poolers. Importantly, we [30] experimentally observed a surprising and counter-intuitive phenomenon for all studied batch senders when maintaining a positive average communication rate across bags. It is advantageous for the sender to use a bag size that is neither too big nor too small to minimize the chances of being caught. Figure 1 is an illustrative example of this phenomenon showing the accuracy of Warden’s pooled detector as a function of number of images sent (the bag size  $B$ ). When pooling evidence from a bag of  $B$  images the pooled detector’s accuracy as a function of  $B$  initially decreases with increasing  $B$ , then levels off, and eventually increases as the Square-Root Law (SRL) [18], [21] inevitably engages since the sender maintains a positive communication rate. The maximal drop in detectability, which we call the *bag gain*, has been observed for all batch senders studied in [30] and for all types of pooled detectors built upon various single-image detectors in the form of rich models as well as convolutional neural networks. It thus appears as a robust phenomenon.

The effect of bag size on security was also previously studied in [27] within the context of Gaussian embedding extended to batch senders. While the authors briefly note what appears to be the bag gain in their experiments, it is not clear how and whether their observation, which was obtained with a single-image source detector, extends to a pooled detector. Indeed, as argued below in this paper and as acknowledged by the authors of [27], to properly assess

The work on this paper was supported by NSF grant No. 2028119.

The authors are with the Department of Electrical and Computer Engineering, Binghamton University, Binghamton, NY, 13902, USA. Email: {edworet1,fridrich}@binghamton.edu

The authors would like to thank Yassine Yousfi for useful discussions and for helping start this research direction.

Copyright (c) 2023 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org.

the performance of batch steganography with pooled detectors, one needs to consider the variability of images within bags, which necessitates adopting a model of cover source diversity, an element missing from [27] as well as [1].<sup>1</sup> Finally, we note that the bag gain did not manifest in previous art [26] because all senders studied in this work embed a variable payload per bag (the rate is maintained only in expectation) based on tags assigned to all images from the cover source computed from an infinitely large bag.

In this paper, we provide an explanation of the bag gain by adopting a model for the soft output of a steganography detector. By suitably simplifying the problem, we become able to analytically study how the bag gain is affected by the detector response, batch sender, cover source, bag size, and communication rate. In a nut shell, the bag gain follows from the square root law, which states that, when the sender maintains a fixed communication rate, the detectability increases with the number of images sent (the bag size). Due to the variability of images across bags, however, the law engages differently in bags of different sizes.<sup>2</sup> The bag gain phenomenon is important for practitioners because the sender can gain on security by spreading the message across multiple images by properly selecting their number (the size of the bag). This gain can be significant and it occurs for bag sizes that can be used in practice.

In the next section, we describe the general setup for batch steganography and pooled steganalysis as considered in this paper. The purpose of Section III is to adopt suitable modeling assumptions that allow us to derive a closed form expression for the performance of Warden’s optimal pooled detector. We also describe the batch sender analyzed in this paper. To capture the diversity of images across bags, in Section IV we adopt a model for the response of a single-image steganography detector on stego images. This model is the key element of our approach as it permits analytic study of the bag gain in Section V, which holds the main bulk of our theoretical results. In particular, we derive a closed-form expression for statistical detectability as a function of the bag size and other parameters describing the cover source and detector response. The derived formulas are contrasted with the performance of a machine learning based pooled detector on real images in Section VI. The model correctly predicts the initial drop in detectability with increasing  $B$ . It also captures experimentally obtained trends in detectability vs. the communication rate (Section V-B3). The model additionally predicts a possible bag loss for bag sizes that are too small, which is experimentally confirmed in real datasets. In Section VII, we extend our analysis to a parametrized family of batch senders to study how the bag gain depends on how strongly the senders adapt the payloads to images from the bag. In Section VIII, we contrast our work with relevant prior art on adaptive bag size. The paper is concluded in Section IX.

Throughout this paper, we use  $\mathcal{N}(\mu, \sigma^2)$  to denote a normal (Gaussian) distribution with mean  $\mu$  and variance

<sup>1</sup>More on the relationship between [27], [1] and our work appears in Section VIII.

<sup>2</sup>Detailed analysis and explanation of the bag gain appears in Section V-A with a summary in plain language presented in the conclusions (Section IX).

$\sigma^2$ . The standard normal tail probability function is denoted  $Q(x) = \int_x^\infty (2\pi)^{-1/2} e^{-z^2/2} dz$ . Symbols  $\mathbb{P}$  and  $\mathbb{E}$  are used for probability and expectation. For a logical statement  $P$ , the indicator function, denoted  $\mathbf{1}_P$ , is equal to 1 when  $P$  is true and 0 when  $P$  is false. The operation of flooring (rounding to the nearest integer  $k \leq x$ ) is denoted  $\lfloor x \rfloor$ .

## II. BATCH STEGANOGRAPHY FORMULATION

Let  $\mathcal{X}$  denote the set of all possible cover images of some fixed size. A cover bag of size  $B$ ,  $\mathbf{X} = (X_0^{(1)}, \dots, X_0^{(B)})$ , is formed by independently selecting  $B$  cover images  $X_0^{(1)}, \dots, X_0^{(B)} \in \mathcal{X}$  according to some probability distribution over  $\mathcal{X}$ . This means that in this paper we do not consider batch senders that select specific covers for embedding since such senders skew the cover source distribution, which would be detectable on its own.

To simplify our analysis and without loss on generality of our conclusions, we will assume that each image from  $\mathcal{X}$  can be embedded at full capacity of  $\log_2 3$  bits per pixel (bpp) with a ternary steganographic scheme. In other words, we assume that images do not contain “wet” pixels [9].

We assume that the steganographer maintains a fixed communication rate  $r \in [0, \log_2 3]$  bpp. This assumption is reasonable as a steganographic channel is likely to be used repetitively in practice. For a fixed positive rate  $r$  expressed in terms of bits per pixel (bpp), the sender will eventually be caught due to the square root law (SRL) [18], [21].

A batch spreading strategy  $S$  is a mapping  $\alpha_{r,S} : \mathcal{X}^B \rightarrow [0, \log_2 3]^B$  that determines the relative payloads (in bpp) embedded in the  $B$  images.<sup>3</sup> When  $r$ ,  $S$ , and  $\mathbf{X}$  are clear from context, we simply write  $\alpha_i \in [0, \log_2 3]$  to denote the  $i$ th component of  $\alpha_{r,S}(\mathbf{X})$ , i.e., the relative payload embedded in the  $i$ th image. The map  $\alpha_{r,S}$  must satisfy the payload constraint  $\sum_{i=1}^B \alpha_i = rB$ . The steganographer produces the  $i$ th stego image  $X_{\alpha_i}^{(i)}$  by embedding cover  $X_0^{(i)}$  with payload of size  $\alpha_i$  bpp using a ternary steganographic scheme.

Next, we provide a general formulation of pooled steganalysis. Given an intercepted bag of  $B$  images  $\mathbf{Y} = (Y^{(1)}, \dots, Y^{(B)})$ , the Warden infers whether steganography is being used by performing the following composite hypothesis test:

$$\begin{aligned} \mathcal{H}_0 : & \quad r = 0 \\ \mathcal{H}_1 : & \quad r > 0. \end{aligned} \tag{1}$$

The Warden “pools” the evidence  $\mathbf{Y}$  together by using a pooled detector (or “pooler”). We assume the Warden’s decision is solely informed by the collection of outputs of a single-image steganography detector, which is a mapping  $d : \mathcal{X} \rightarrow \mathbb{R}$  that assigns to each image a scalar referred to as the soft output (or response) of the detector. Formally, the Warden’s pooler is of the form  $\pi : \mathbb{R}^B \rightarrow \mathbb{R}$ , and she infers whether the sender uses steganography by computing  $d(Y^{(i)})$  for all  $i = 1, \dots, B$  and comparing  $\pi(d(Y^{(1)}), \dots, d(Y^{(B)}))$  against a threshold determined by some application-dependent requirements, such as controlling the false alarm.

<sup>3</sup>Notice that the mapping is deterministic as we are not considering randomized spreading strategies in this paper.

In the next two sections, we simplify the formulation above in order to study the bag gain phenomenon analytically. Our approach is *detector-centric* in the sense that we

- 1) impose statistical models on the response of the detector  $d$  on cover and stego images and let all actors share information (next section)
- 2) model the diversity of bags with a suitably simplified statistical model of the so-called detector response curves that express the dependence of the detector output on message length (Section IV).

### III. MODELING ASSUMPTIONS

This paper's goal is to analytically capture and intuitively explain the experimentally observed bag gain phenomenon. This necessitates a rather significant simplification of the setup described in the previous section in terms of what knowledge is available to all actors and in terms of modeling assumptions to facilitate an analytically tractable analysis. To this end, we introduce the concept of acquisition oracle and make specific assumptions about statistical properties of a single-image detector when applied to cover and stego images. We also introduce the batch sender studied in this paper.

The act of taking an image with a digital camera introduces randomness into the image due to numerous acquisition noise sources, such as the shot (photonic) noise, the readout noise, and thermal noise [15]. Thus, taking multiple images of the exact same scene with the same camera would produce slightly different images that follow a statistical distribution, which we call in this paper an *acquisition oracle*, a concept that found many uses in steganography in the past [2], [10], [11], [28], [29]. Given a collection of cover images indexed by  $i = 1, \dots, B$ , we consider the specific cover image  $X_0^{(i)}$  used by the sender as a sample from the oracle.<sup>4</sup> This oracle will provide us with the means to narrow down the distribution of  $d(Y^{(i)})$  under both hypotheses. As will become apparent in the next section, in this paper we will only need to make assumptions on the distribution of detector outputs on the realizations of the oracle, avoiding thus the potentially complex task of modeling the oracle itself.

#### A. Gaussianity and local shift hypothesis

First, we take advantage of the fact that, for each  $i$ , the distribution of the  $i$ th cover image  $X_0^{(i)}$  is concentrated on a small subset of  $\mathcal{X}$  (multiple images of the same  $i$ th scene taken with the same camera differ only slightly). Since differentiable non-linear functions are approximately linear on sufficiently small neighborhoods, we can employ the central limit theorem (CLT) so that<sup>5</sup>

$$d(X_0^{(i)}) \sim \mathcal{N}(\mu_i, \sigma_i^2), \quad (2)$$

where  $\mu_i$  and  $\sigma_i^2$  are the expected value and variance of  $d$  on cover images generated by the acquisition oracle for the  $i$ th scene. Since stego schemes try to preserve statistical

properties of  $X_0^{(i)}$ , the embedding process will also preserve the concentration. Therefore, by the same argument we assume that  $d(X_{\alpha_i}^{(i)})$  is also Gaussian<sup>6</sup>

$$d(X_{\alpha_i}^{(i)}) \sim \mathcal{N}(\mu_i + s_i(\alpha_i), \sigma_i^2) \quad (3)$$

with an additional assumption that only the mean is affected by embedding but not the variance. This *local* shift hypothesis is a much weaker assumption than the shift hypothesis [26] about the *global* distribution of detector response which is not satisfied for modern steganalyzers in the form of rich models and CNNs (see Sec. 3.2 in [30]).

Technically, the variance of  $d(X_{\alpha_i}^{(i)})$  also depends on  $\alpha_i$  because of the added randomness in the form of the stego key selection and the message itself. We do not consider this dependence in order to further simplify the modeling and also because the acquisition noise dominates the statistical spread because it is stronger than the stego noise.

Finally, to avoid modeling the distribution of the variances  $\sigma_i^2$  across images from  $\mathcal{X}$  and the oracle itself, we assume all variances are the same across scenes  $\sigma_i^2 = \sigma^2$ .

#### B. Uniformity of response increase

The response curve (RC) for image  $X_0^{(i)}$  and detector  $d$  is the function  $\varrho_i : [0, \log_2 3] \rightarrow \mathbb{R}$  defined by

$$\varrho_i(\alpha) = \mathbb{E}[d(X_{\alpha}^{(i)}) | X_0^{(i)}]. \quad (4)$$

Given the payload size  $\alpha$  and a fixed cover  $X_0^{(i)}$ ,  $\varrho_i(\alpha)$  is the expected value of the response  $d(X_{\alpha}^{(i)})$  when embedding  $X_0^{(i)}$  with random messages and stego keys.

Since the detector is trained to be sensitive to embedding changes but not acquisition noise, we assume the expected increase in detector response is uniform across all possible acquisitions

$$\varrho_i(\alpha) - \varrho_i(0) = s_i(\alpha) \quad (5)$$

for all realizations of  $X_0^{(i)}$ . This assumption allows us to compute the expected shift  $s_i(\alpha)$  from a specific cover image, which simplifies analysis and practical implementations.

#### C. Warden's test

Equipped with a single-image detector  $d$  that adheres to the assumptions above, the Warden's hypothesis test (1) becomes:

$$\begin{aligned} \mathcal{H}_0 : & \quad d(Y^{(i)}) \sim \mathcal{N}(\mu_i, \sigma^2) \quad \text{for all } i \\ \mathcal{H}_1 : & \quad d(Y^{(i)}) \sim \mathcal{N}(\mu_i + s_i(\alpha_i), \sigma^2) \quad \text{for all } i, \end{aligned} \quad (6)$$

where  $Y^{(i)}$  are the images from a bag under inspection by the Warden and  $\alpha_i$  is the payload residing in the  $i$ th image.

Assuming the parameters of the distributions in the hypothesis test (6) are known to the Warden, the test becomes simple and the Warden's most powerful pooled detector is the likelihood ratio test. The detectability of steganography in a single bag is determined by the deflection coefficient

$$\Delta^2(\mathbf{X}) = \sum_{i=1}^B \frac{s_i^2(\alpha_i)}{\sigma^2} = \sum_{i=1}^B \frac{(\varrho_i(\alpha_i) - \varrho_i(0))^2}{\sigma^2}, \quad (7)$$

<sup>4</sup>Note that the oracle is a conditional distribution describing the distribution of images (acquisitions) of a specific scene.

<sup>5</sup>Modern single-image detectors  $d$  are often neural networks with differentiable structure.

<sup>6</sup>The random variable  $X_{\alpha_i}^{(i)}$  is generated by 1) sampling  $X_0^{(i)}$  from the oracle and 2) embedding a random message with a random stego key.

where  $s_i(\alpha_i)$  can be computed via  $\varrho_i(\alpha_i) - \varrho_i(0)$  given any oracle realization  $X_0^{(i)}$ .

#### D. Minimum deflection sender

As a batch sender for our study, we selected the detector-informed Minimum Deflection Sender (MDS) introduced in [30] because it is the most amenable to analysis within the context of a statistical model of the detector. As will be argued in Section VII, the bag gain generally manifests for batch senders that minimize the risk of being detected by assigning larger payloads to images that are difficult to steganalyze and smaller payloads to images in which the embedding is more detectable. In particular, the bag gain has also been observed for the detector-agnostic Image Merging Sender (IMS) [26] and detector-aware Shift Limited Sender (SLS) [30].

The MDS makes use of a single-image detector, which we will assume is the same as the one used by the Warden. Given a bag of images  $\mathbf{X}$ , the MDS selects payloads  $\alpha_i$  that minimize the deflection (7). Formally,  $\alpha_i$  are found by solving the following optimization problem

$$\begin{aligned} & \text{minimize } \Delta^2(\mathbf{X}), \\ & \text{s.t. } \sum_{i=1}^B \alpha_i = rB, \alpha_i \in [0, \log_2 3] \forall i, \end{aligned} \quad (8)$$

where  $r \in [0, \log_2 3]$  is a chosen embedding rate in bpp. A general solution is given in Appendix B.

Granting the Warden and the MDS access to the same detector  $d$  makes the MDS the optimal batch sender—it minimizes the power of the Warden’s most powerful detector.

#### E. Discussion

Our setup assumes the actors are omniscient. Among other things, the Warden knows the steganographic method used by the sender, the payloads  $\alpha_i$  possibly embedded in each image, the communication rate  $r$ , and the bag size  $B$ . Moreover, the sender and the Warden share the same single-image detector. While it is certainly of interest to study more relaxed setups and perhaps even probabilistic strategies within game theory, such scenarios would require adopting and justifying additional models on how accurately the Warden can estimate the payloads  $\alpha_i$ , on the nature of the mismatch between the detectors, etc. The fact that our conclusions regarding the bag gain based on the simplified setup do capture trends observed in real-life situations testify to their relevance.

Having said this, we wish to point out to the reader that the bag gain has been observed in experiments under much more relaxed conditions, including different pooling strategies, mismatched and qualitatively different detectors built using various machine-learning paradigms, and when the Warden needs to estimate the embedded payloads from the images at hand. The reader is advised to inspect Section 7 in [30] for more details.

## IV. RESPONSE CURVE MODEL

In order to analyze the trends of detectability w.r.t. the bag size  $B$  and possibly the communication rate  $r$ , we must somehow obtain a model of  $\Delta^2(\mathbf{X})$  over bags since  $\mathbf{X}$  has an underlying distribution. We must be careful with our modeling assumptions to preserve the essential complexities of Eq. (1) so that the bag gain can properly manifest. Due to the form of the deflection  $\Delta^2(\mathbf{X})$  in Eq. (7), it is sufficient to model the response curves across images, which is easier than modeling natural images and also keeps a tighter connection between the model and practice. In particular, we make the following two assumptions about response curves.

#### A. Linear response curves

We first assume the response curves are linear in payload

$$\varrho_i(\alpha_i) - \varrho_i(0) = b_i \alpha_i, \quad (9)$$

where  $\alpha_i \in [0, \log_2 3]$  and  $b_i \in [0, \infty)$  is the slope of the linear response curve. This significantly simplifies the problem, permitting a closed-form expression for the payloads  $\alpha_i$  embedded by the MDS and its extension in Section VII. Even though the response curves of typical detectors built with machine learning are not linear (see, e.g., Figure 3 in [30]), they are approximately linear when  $\varrho_i(\alpha_i) - \varrho_i(0)$  is small.

#### B. Binomial model for slopes

Arguably, if all images from the cover source had similar response curves, the MDS would spread payload nearly uniformly, at which point the detectability would need to increase with  $B$  from the beginning due to the SRL. The reason for the bag gain is source diversity and the fact that the counts of images that contain very small payloads and those that are embedded nearly fully fluctuate across bags. Thus, in order to simplify the modeling but preserve the essence we adopted a two-valued range for the response curve slopes  $b_i$ :  $\mathbb{P}(b_i = \varepsilon) = p$  and  $\mathbb{P}(b_i = 1) = 1 - p$  where  $0 < \varepsilon \ll 1$  and  $p \in [0, 1]$ . Let  $C_\varepsilon$  denote the number of response curves with slope  $\varepsilon$  in a bag of size  $B$ . Assuming the images are drawn randomly from the cover source,  $C_\varepsilon$  follows a binomial distribution on  $\{0, 1, \dots, B\}$ .

The linear model of response curves and the binomial model of slopes was adopted to simplify the problem and permit the subsequent analysis. In real life, the batch sender can of course use real response curves and minimize (7) numerically. In fact, this is exactly how the sender was implemented in our previous work [30] that lead to the discovery of the bag gain.

It is easy to show that if all  $B$  images have uniform slope  $b$ , the deflection  $\sum b^2 \alpha_i^2$  is minimal when all images receive uniform payload  $\alpha_i = \alpha$ . In a bag of two images with different slopes, they both start receiving non-zero payload when embedding a message of any length. More generally via a water filling algorithm (see Appendix B), with increasing rate  $r$  all images in the bag start receiving payload until the ones with slope  $\varepsilon$  saturate at  $\log_2 3$ . From there, the images with slope 1 absorb the remaining payload.

### C. Pooled detector performance measure

The deflection coefficient  $\Delta^2(\mathbf{X})$  (7), which depends on  $\varepsilon, r, B$ , and  $C_\varepsilon$ , informs us about the performance of the likelihood ratio detector in a specific bag of images. For fixed  $\varepsilon, r, B$ , the Receiver Operating Characteristic curve (ROC) of the pooled detector expressing the probability of correct stego bag detection  $P_D$  as a function of the probability of false alarm  $P_{FA}$  is the expectation over bags

$$P_D(P_{FA}) = \mathbb{E}[Q(Q^{-1}(P_{FA}) - \Delta(\mathbf{X}))] \quad (10)$$

$$= \sum_{k=0}^{\infty} \frac{(-1)^k c_k}{k!} Q^{(k)}(Q^{-1}(P_{FA}) - \mathbb{E}[\Delta(\mathbf{X})]),$$

where  $c_k$  is the  $k$ th central moment of  $\Delta(\mathbf{X}) \triangleq \sqrt{\Delta^2(\mathbf{X})}$  as shown in Appendix A. Keeping only the terms up to  $k = 2$  in the sum provides a rather accurate approximation for typical values of our modeling parameters (note that  $c_1 = 0$ ).

In this paper, our reasoning is based on the expectation of the deflection coefficient because it is significantly easier to analyze than the ROC (10). While the expected deflection informs us about the ROC over bags *indirectly* (as seen from (10)), many qualitative properties observed for the expected deflection do propagate to common scalar ROC measures, such as the weighted Area Under the Curve (wAUC) [6].

## V. EXPLAINING THE BAG GAIN

In this section, we explain the performance trends using the binomial linear model for response curves. We begin by simply assuming that images can hold an arbitrarily large amount of payload. As we progress through this section, we will incorporate more realistic constraints in order to capture which pieces of the model are responsible for certain phenomena we observe in practice.

### A. Unbounded embedding capacity

First, we analyze the case of unbounded embedding capacity for all images from the bag. We believe it is useful to start with this case as it 1) clearly captures important trends in the small bag regime, 2) is analytically tractable, and 3) serves to build the reader's intuition as to why a bag gain should occur in the first place. Studying the unbounded case will also help underscore the impact of finite embedding capacity on the observed trends later seen in Section V-B.

Based on Eq. (44) in Appendix B, the MDS payloads for the unbounded case are given, for all  $i$ , by

$$\alpha_i = \frac{rB}{b_i^2 \sum_{k=1}^B \frac{1}{b_k^2}} = \frac{rB\varepsilon^2}{b_i^2(C_\varepsilon + (B - C_\varepsilon)\varepsilon^2)}, \quad (11)$$

since

$$\sum_{k=1}^B \frac{1}{b_k^2} = C_\varepsilon \varepsilon^{-2} + (B - C_\varepsilon). \quad (12)$$

Utilizing (11) and (12), the deflection simplifies to

$$\Delta^2(\mathbf{X}) = \frac{1}{\sigma^2} \sum_{i=1}^B b_i^2 \alpha_i^2 = \frac{r^2 B^2 \varepsilon^2}{\sigma^2 (C_\varepsilon + (B - C_\varepsilon)\varepsilon^2)}. \quad (13)$$

In this case, the expected deflection becomes

$$\mathbb{E}[\Delta^2(\mathbf{X})] = \frac{r^2 B^2 \varepsilon^2}{\sigma^2} \sum_{k=0}^B \frac{\binom{B}{k} p^k (1-p)^{B-k}}{k + (B-k)\varepsilon^2}, \quad (14)$$

which can be further simplified using Stirling's formula (see, e.g., page 147 in [8]) as  $B \rightarrow \infty$

$$\begin{aligned} \binom{B}{pB} &\sim 2^{BH_2(p)} \\ \Rightarrow \binom{B}{pB} p^{pB} (1-p)^{(1-p)B} &\sim 2^{BH_2(p)} \times 2^{-BH_2(p)} = 1 \\ \Rightarrow \mathbb{E}[\Delta^2(\mathbf{X})] &\sim \frac{r^2 \varepsilon^2 B}{\sigma^2 (p + \varepsilon^2(1-p))}. \end{aligned} \quad (15)$$

Here,  $H_2$  is the binary entropy function, and  $\sim$  means the ratio of both sides tends to 1 as  $B \rightarrow \infty$ .

Figure 2 shows the expected deflection  $\mathbb{E}[\Delta^2(\mathbf{X})]$  as a function of the bag size  $B$  with the dashed lines drawn to show asymptotic trends 15. The figure also shows wAUC of Eq. (10) as a function of  $B$ . For small  $p$ , the detectability initially grows due to the SRL because the bags are small and most do not contain any images with slope  $\varepsilon$ . The growth is steep because it is driven primarily due to payload embedded in images with slope 1. As the bag size increases, however, the detectability starts dropping since the bags are more likely to contain images with small slopes which absorb most of the payload with only a slight contribution to the deflection. The deflection eventually levels off and then linearly increases. This time, the growth is less steep because of the presence of images with small slope  $\varepsilon$ . Thus, the existence of the local maximum and global minimum of expected deflection is fundamentally a consequence of the *SRL switching its growth rate*.

As depicted in Figure 2, the unbounded capacity model predicts two critical bag sizes that depend primarily on  $p$  and  $\varepsilon$ . One is associated with a local maximum,  $B_{\max}$ , while the other,  $B_{\min}$ , corresponds to minimal expected deflection. We do not talk about these critical bag sizes as corresponding to bag loss and bag gain *yet* because we define these concepts for the more realistic bounded capacity case using an easily interpretable performance measure (wAUC) in the next section. The closed form for the expected deflection as a function of bag size allows us to study the critical points and obtain insight into the conditions under which the local maximum and the minimum can occur and how they depend on  $\varepsilon$  and  $p$ . Figure 2 tells us that we can then implicitly (but indirectly) draw conclusions about wAUC since the relationships closely transfer as visually portrayed.

Since our model is only defined for positive integers  $B \geq 1$  (actual bag sizes), we begin by simplifying the expression in Eq. (14) by using Eq. (15) (the dominant term in the large bag regime) along with the  $k = 0$  term (the dominant term in the small bag regime when  $\varepsilon$  is small) :

$$\mathbb{E}[\Delta^2(\mathbf{X})] \doteq \frac{r^2 B}{\sigma^2} \left( (1-p)^B + \frac{\varepsilon^2}{p + \varepsilon^2(1-p)} \right). \quad (16)$$

Notice that Eq. (16) can be defined on the real numbers  $B \in \mathbb{R}$ . Using mild simplifying assumptions, we can derive closed

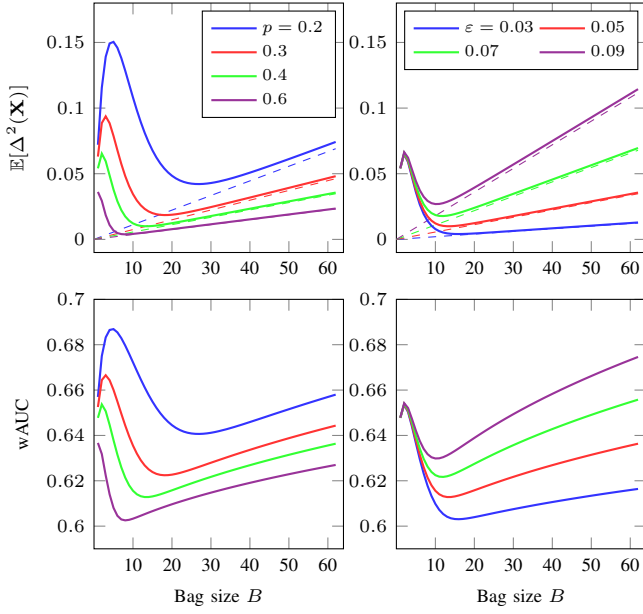


Figure 2. Unbounded capacity model of pooled detector performance. Top row is  $\mathbb{E}[\Delta^2(\mathbf{X})]$  as a function of  $B$  (solid) with the line (dashed) drawn to show asymptotic trends (15). Bottom row is wAUC of Eq. (10) as a function of  $B$ . Left column shows trends w.r.t.  $p$  ( $\varepsilon = 0.05$  fixed) and right column shows trends w.r.t.  $\varepsilon$  ( $p = 0.4$  fixed). We have  $r = 0.3$  and  $\sigma^2 = 1$  fixed.

form approximations for both critical bag sizes. Specifically, using Eq. (16) and the fact that  $(1-p)^B = e^{B \ln(1-p)}$ , we can approximate the optima by finding solutions to

$$\frac{\partial}{\partial B} \mathbb{E}[\Delta^2(\mathbf{X})] \doteq \left( e^{B \ln(1-p)} + \frac{\varepsilon^2}{p + \varepsilon^2(1-p)} \right) \quad (17)$$

$$+ B \ln(1-p) e^{B \ln(1-p)} = 0. \quad (18)$$

Since  $\ln(1-p) < 0$ , for small  $B$  the term proportional to  $\varepsilon^2$  is small compared to the other two terms. Setting  $\frac{\varepsilon^2}{(p + \varepsilon^2(1-p))} \approx 0$ , we obtain an approximate formula for the first critical bag size corresponding to the local maximum<sup>7</sup>

$$0 = e^{B \ln(1-p)} (1 + B \ln(1-p))$$

$$\Leftrightarrow B_{\max} \doteq \frac{-1}{\ln(1-p)}. \quad (19)$$

For larger bag sizes, the term proportional to  $\varepsilon^2$  cannot be ignored. We rearrange the terms and take log of both sides (keep in mind that  $\ln(1-p) < 0$ )

$$\frac{\varepsilon^2}{(p + \varepsilon^2(1-p))} = -e^{B \ln(1-p)} (1 + B \ln(1-p))$$

$$\Leftrightarrow B \ln(1-p) = \ln \left( \frac{\varepsilon^2}{p + \varepsilon^2(1-p)} \right) \quad (20)$$

$$- \ln(-1 - B \ln(1-p)). \quad (21)$$

Since the second term on the r.h.s. of this equation is small with respect to the l.h.s., we obtain a first order approximation

<sup>7</sup>The fact that  $B_{\max}$  corresponds to a local maximum can be verified by computing the second derivative.

for the second critical bag size<sup>8</sup>

$$B_{\min} \doteq \frac{\ln \left( \frac{\varepsilon^2}{p + \varepsilon^2(1-p)} \right)}{\ln(1-p)}. \quad (22)$$

From (19), we can deduce that the initial growth associated with the local maximum ceases to manifest with sufficiently large prior probability  $p$  of images with small slopes. In particular,  $B_{\max} < 1$  for  $p \gtrsim 0.63$  in approximate agreement with Figure 2 when working with the exact expected deflection. Additionally, Eq. (22) encapsulates how  $B_{\min}$  depends on  $p$  and  $\varepsilon$  (it increases as  $\varepsilon$  or  $p$  decrease). This makes intuitive sense as smaller  $\varepsilon$  means the images can hold larger payload, making the SRL take longer to finish switching its growth rate. Similarly, with a smaller fraction  $p$  of such images, it takes larger bags to see their effect on detectability.

### B. Bounded embedding capacity

We now show the effect of bounding the embedding capacity to  $A = \log_2 3$  bpp and also formally define the bag loss and bag gain. Images with  $b_i = \varepsilon$  achieve embedding capacity  $\alpha_i = A$  when (c.f. Eq. (11))

$$\frac{rB}{C_\varepsilon + (B - C_\varepsilon)\varepsilon^2} \geq A, \quad (23)$$

which holds iff

$$T := \frac{r/A - \varepsilon^2}{1 - \varepsilon^2} B \geq C_\varepsilon. \quad (24)$$

If  $T < C_\varepsilon$ , then  $\Delta^2(\mathbf{X})$  is given by Eq. (13). However, if  $T \geq C_\varepsilon$ , we have

$$\alpha_i = \begin{cases} \frac{rB - AC_\varepsilon}{B - C_\varepsilon} & b_i = 1 \\ A & b_i = \varepsilon \end{cases} \quad (25)$$

and so

$$\Delta^2(\mathbf{X}) = \frac{C_\varepsilon \varepsilon^2 A^2}{\sigma^2} + \frac{(rB - C_\varepsilon A)^2}{\sigma^2 (B - C_\varepsilon)}. \quad (26)$$

Thus, we have in expectation

$$\mathbb{E}[\Delta^2(\mathbf{X})] = \mathbb{E}[\Delta^2(\mathbf{X}) \mathbf{1}_{T < C_\varepsilon}] + \mathbb{E}[\Delta^2(\mathbf{X}) \mathbf{1}_{T \geq C_\varepsilon}]$$

$$= \frac{1}{\sigma^2} \left[ r^2 B^2 \varepsilon^2 \sum_{k=[T]+1}^B \frac{\binom{B}{k} p^k (1-p)^{B-k}}{k + (B-k)\varepsilon^2} \right.$$

$$+ \sum_{k=0}^{[T]} \binom{B}{k} p^k (1-p)^{B-k}$$

$$\left. \times \left( k \varepsilon^2 A^2 + \frac{(rB - kA)^2}{B - k} \right) \right]. \quad (27)$$

In Figure 3, we show the wAUC of Eq. (10) (instead of expected deflection) for various combinations of  $\varepsilon, r, p$  since we intend to contrast the performance of the model with real life detectors.

<sup>8</sup>A more precise argument can be made here based on iterative root finding for the equation  $B = f(B)$  by showing that  $|f'(B)| < 1$  for convergence.

1) *Bag gain and bag loss*: While the exact trend of wAUC w.r.t.  $B$  depends on  $\varepsilon$ ,  $r$ , and  $p$ , one can roughly say that (ignoring for now the small oscillations commented upon in the next section): 1) wAUC can either grow right from  $B = 1$ , or 2) grow, reach a local maximum, decrease, reach a global minimum (bag gain), and then increase, or 3) exhibit a global minimum without the initial increase. Fundamentally, the local maximum and the global minimum of wAUC are due to the varying statistical makeup of small bags as already commented for the unbounded capacity case. Eventually, for large enough  $B$  wAUC will approach 1. How fast this happens depends on whether large enough bags contain enough images with small slopes to avoid embedding substantial payload in images with a large slope. This occurs approximately when  $p \log_2 3 > r$ , at which point wAUC approaches 1 only very slowly, depending on the value of  $\varepsilon$ . This is why the global minimum appears quite shallow for some combinations of the parameters.

Formally, we define the bag gain  $\gamma$  as the maximum decrease in a chosen detectability measure the batch sender can enjoy by bagging. Since we use wAUC,

$$\gamma = \max_{B \geq 1} [\text{wAUC}(1) - \text{wAUC}(B)], \quad (28)$$

where  $\text{wAUC}(B)$  is the wAUC of the pooled detector on bags of size  $B$ . Notice that the bag gain can be observed for most combinations of the parameters in Figure 3 but disappears for large enough rates and for larger  $\varepsilon$ . This is intuitively correct as larger rates force the detectability to grow faster as do larger values of  $\varepsilon$ .

Besides the global minimum corresponding to the bag gain, wAUC as a function of  $B$  may exhibit a local maximum for small bags (for  $p \leq 0.3$  in the figure). When the bag gain is positive ( $\gamma > 0$ ), we define bag loss as

$$\nu = \max_{B_{\text{Gain}} > B \geq 1} [\text{wAUC}(B) - \text{wAUC}(1)], \quad (29)$$

where  $B_{\text{Gain}}$  is the bag size corresponding to the bag gain.<sup>9</sup> In words, bag loss is the increase in detectability when the sender selects the worst bag size instead of the optimal  $B_{\text{Gain}}$ . Based on our definition, bag loss is not defined if there is no positive bag gain. Similar to the bag gain, bag loss may not manifest for certain combinations of the parameters.

2) *Local oscillations*: As shown in Figure 3, the wAUC experiences a transient oscillating / periodic behavior for smaller bag sizes, which can be explained by analyzing expected deflection. The oscillations appear when considering images with bounded capacity and are ultimately due to the quantization of  $T$  when computing the bounds for the sums in Eq. (27). In particular, since  $\varepsilon^2$  is small,  $T \approx rB/A$  which implies  $\lfloor T \rfloor$  increments whenever  $B \approx Ak/r$  for some positive integer  $k$ . In other words,  $\lfloor T \rfloor$  is fixed for intervals of length  $A/r$ . For example, for  $r = 0.3$  we have  $A/r \approx 5$  which is approximately the period shown in the corresponding plot in Figure 3. Within each interval,  $\mathbb{E}[\Delta^2(\mathbf{X})]$  (and wAUC) changes in a continuous manner and may contain local optima due to the upper sum  $\mathbb{E}[\Delta^2(\mathbf{X})\mathbf{1}_{T < C_\varepsilon}]$ .

3) *Trends w.r.t rate*: In the unbounded case, we see that the expected deflection is linearly proportional to  $r^2$  (15). However, in the bounded capacity case, the rate has a non-trivial affect on the performance curves (in terms of wAUC) as seen in Figure 3 and, in particular, the location of  $B_{\text{Loss}}$  and  $B_{\text{Gain}}$ . For example, as  $r$  increases we see that  $B_{\text{Gain}}$  decreases for  $\varepsilon = 0.06$  and  $p = 0.2$ , but  $B_{\text{Gain}}$  increases for  $\varepsilon = 0.02$  and  $p = 0.4$ . Note that if  $T < 1$ , then approximately  $rB < A$  which makes Eq. (27) degenerate to the unbounded model Eq. (14).

## VI. OBSERVING TRENDS IN REAL IMAGES

In this section, we contrast the trends in detectability w.r.t. bag size from experiments with real images and detectors with those obtained from the model. We measure the performance with wAUC. First, we describe our experimental setup, including the dataset and a single-image detector used by some batch senders and for pooled steganalysis. As mentioned in [30], the embedding algorithm (whether cost-based or model-based) does not have a significant effect on the bag gain manifesting, so we limit our experiments to the cost-based HILL [23]. All experiments were done on the image dataset ALASKA II [6] developed as in [6] without the final JPEG compression step.<sup>10</sup> We consider two disjoint subsets of ALASKA II images denoted split1 and split2, containing 25,000 images each. Split1 is used to train the shared single-image detector and Warden’s pooled detector while split2 is used to assess the performance of batch senders.

The detector-aware senders use a single-image detector  $d$  in the form of an SRNet [4] pre-trained on ImageNet with the binary task of steganalyzing J-UNIWARD [12] (the so-called JIN pre-training exactly as described in [5]). The refinement to detect HILL was done on a diverse stego source created using split1 with relative payloads randomly drawn from the uniform distribution on the set of relative payloads

$$\mathcal{P} = \{0.05, 0.1, 0.2, \dots, 1.4, 1.5\}. \quad (30)$$

In particular, split1 was partitioned into further subsets of 22k, 1k, and 2k images for training, validation, and testing, respectively. The detector-aware senders use the logit as the detector’s response.

The Warden is given the sender’s detector  $d$  for steganalysis. She is also assumed clairvoyant and given the knowledge of the payloads  $\alpha_i$ . The reader is referred to [30] for a comprehensive analysis of the situation when the Warden estimates  $\alpha_i$  from the images at hand and when she trains her own single-image detector that is possibly different as well as trained on a different dataset from the same source. In particular, as shown in this prior art, the trends of detectability vs. bag size appear to be robust and unaffected by Warden’s choices.

Three batch senders are tested: the Image Merging Sender (IMS) and the detector-aware Shift Limited Sender (SLS) and MDS. The IMS treats each bag as one big image and lets the given stego algorithm decide what payload chunk each

<sup>9</sup> $B_{\text{Loss}}$  will denote the bag size corresponding to the bag loss.

<sup>10</sup>The authors note that the bag gain was observed on other datasets, such as BOSSbase [7] and BOWS2 [3] (not shown in this paper).

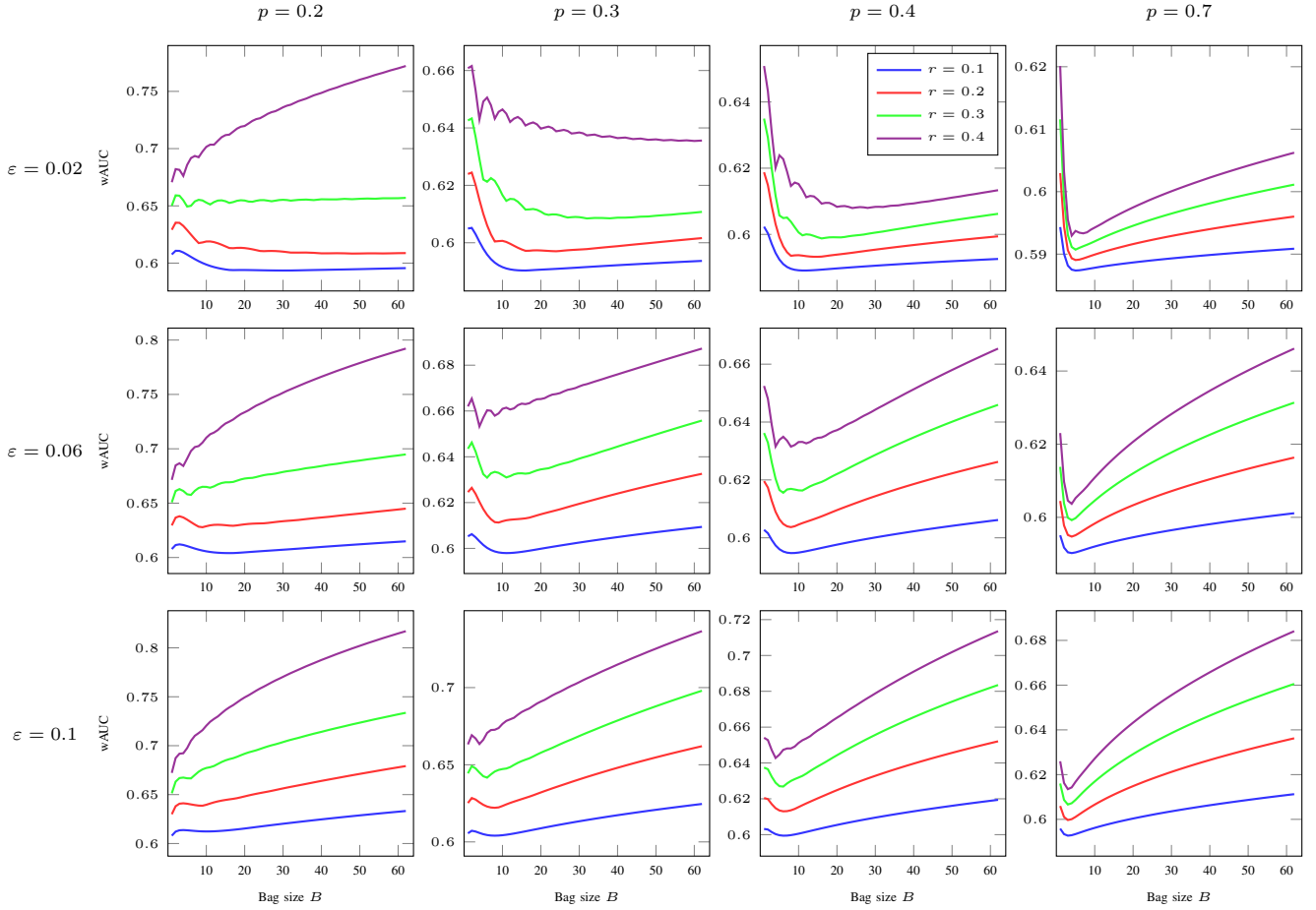


Figure 3. Bounded capacity model of optimal pooled detector performance (wAUC) as a function of  $B$  for various combinations of  $\varepsilon, r, p$ . Rows correspond to fixed  $\varepsilon$ , columns correspond to fixed  $p$ , and colors correspond to fixed  $r$ .

image will hold. The SLS finds the payloads by requiring that the embedding induces the same shift in the detector response. The MDS, which is described in Section III-D, was implemented using a projected gradient descent method to find optimal payloads since response curves for real images are non-linear. We refer the reader to the original publication for more details [30]. We did not include the batch sender proposed in [27] because it is equivalent to the IMS with an embedding scheme adjusted as in Gaussian embedding.

The optimal pooled detector described in Section III-C was used to analytically study and explain the bag gain trends; however, such a pooler is infeasible in practice due to the difficulty of estimating the parameters of the distributions in (6). Thus, all experiments on real images use the LRT pooler,  $\pi_{\text{LRT}}$ , as thoroughly studied in [30]. The Warden tests whether the detector output for the  $i$ th image of the bag is consistent with the distribution of the detector  $f_{\alpha_i}$  on stego images all embedded with the same relative payload  $\alpha_i$ :

$$\begin{aligned} \mathcal{H}_0 &: d(Y^{(i)}) \sim f_0 \quad \text{for all } i \\ \mathcal{H}_1 &: d(Y^{(i)}) \sim f_{\alpha_i} \quad \text{for all } i \end{aligned} \quad (31)$$

with the optimal detector being the log-likelihood ratio

$$\pi_{\text{LRT}}(\mathbf{Y}) = \sum_{i=1}^B \log \frac{f_{\alpha_i}(d(Y^{(i)}))}{f_0(d(Y^{(i)}))}. \quad (32)$$

The distributions  $f_{\alpha_i}$  are estimated empirically using the test set of split1.<sup>11</sup> Both spreading and pooling is done on split2.

We note that [30] investigated three other pooled detectors, including situations when the Warden trained the detector on a different dataset and/or used a different neural architecture or even a qualitatively different detector, such as a rich model. The bag gain was generally observed under all circumstances. For a comprehensive look at bag gain trends across poolers in general, we refer the reader to [30].

#### A. Trends seen in ALASKA II

Our focus is on trends of detectability w.r.t. bag size  $B$  and rate  $r$  for multiple batch senders. Figure 4 shows the detection performance of the LRT pooler  $\pi_{\text{LRT}}$ . For each fixed  $B, r$ , and sender, we independently form 2000 bags sampled without replacement from split2. The wAUC is computed from the ROC formed by the 2000 samples of bags.

First, notice that all senders exhibit a bag gain, including the detector-agnostic IMS. The steganographer can use the bag gain to decrease the statistical detectability by up to  $\sim 0.15$  in pooled detector performance, which can significantly benefit the steganographer in practice. Second, the initial decrease in detectability engages quickly so even using bags of size

<sup>11</sup>Using SciPy's gaussian\_kde function



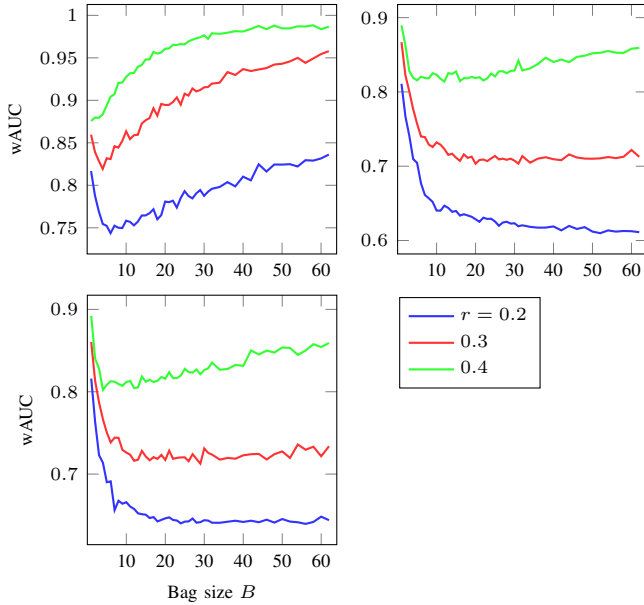


Figure 4. Trends in the performance of  $\pi_{\text{LRT}}$  across batch senders for ALASKA II (top left IMS, top right SLS, bottom left MDS). For the lower payloads of SLS and MDS, the SRL requires a much larger bag size to take effect.

5, e.g., as opposed to using a single-image is significantly advantageous for the steganographer.

Despite the differences between response curves under the binomial model and real image response curves, the trends predicted by our model and shown in Figure 3 provide valuable insight. In particular, the model correctly predicts that for large enough payloads the bag gain disappears. Furthermore, the optimal bag size  $B_{\text{Gain}}$  increases with decreased rate  $r$  except for the smallest value of  $\varepsilon$  (cover source with images with basically flat response curves). Our model additionally predicts that this increase is smaller in cover sources with fewer hard-to-steganalyze images (smaller  $p$ ).

One of the clearest differences between IMS and the two detector-aware senders that can be seen in Figure 4 is that the SRL engages a lot sooner for IMS. The main contributing factor is that the two detector-aware senders are more aggressive in utilizing difficult images by embedding them with larger payloads because they are aware of the impact on detectability. Batch senders that are even less aggressive than IMS will eventually not exhibit the bag gain. In the extreme case of a batch sender that assigns the same payloads to all images, the detectability will monotonically increase as per the large bag regime’s SRL. In Section VII, we will explain this behavior from a model by introducing a family of batch senders parametrized by a scalar parameter (the Hölder sender) that encompass the uniform sender, the SLS, and MDS.

Finally, as seen in Figure 3 for some combinations of  $\varepsilon$ ,  $p$ , and  $r$  our model predicts oscillations in wAUC for small bag sizes and an initial bag loss (local maximum in wAUC) for very small bag sizes. While these higher-order effects were not observed in our experiments on ALASKA II, in the next section we demonstrate that they are real phenomena that can manifest in other datasets with the right diversity of images.

## B. Bimodal ALASKA II

As commented on in the previous section, our binomial model of slopes predicts that, for small bag sizes and certain combinations of  $\varepsilon$ ,  $p$ , and  $r$ , wAUC should exhibit a local maximum, the bag loss, and oscillations that decay with larger bag sizes. Such higher-order effects are neither seen in our experiments nor in the prior art [30] because the real distribution of response curves in images from ALASKA II is not close enough to the binomial model of slopes. Of course, it does not mean that bag loss cannot occur in other datasets.

In order to investigate whether these phenomena can manifest for real images, we construct multiple versions of approximately “bimodal” ALASKA II consisting of two groups of images: 1) easy-to-steganalyze images with steep response curves and 2) hard-to-steganalyze images with almost flat response curves. Such approximately bimodal distribution can realistically occur, for example, in a landscape photographer’s portfolio when the majority of the source is low ISO images, which would be the case of images taken during daylight, while the remainder is high ISO images taken during the night (astrophotography).

We propose the following stochastic procedure based on rejection sampling to enforce a distribution of slopes on ALASKA II that more closely matches our model. This will also allow us to parameterize the dataset by  $p$ , a source diversity parameter, so we can feasibly observe trends across sources with a varying proportion of easy-to-steganalyze and hard-to-steganalyze images.

First, we perform what we call “ $\varepsilon/M$  binning” on ALASKA II. Given four non-negative constants  $\ell_\varepsilon \leq u_\varepsilon \leq \ell_M \leq u_M$ , we say image  $X$  has an  $\varepsilon$ -type RC  $\varrho_X$  if for all  $\alpha \in \mathcal{P}$ ,  $\ell_\varepsilon \alpha \leq \varrho_X(\alpha) - \varrho_X(0) \leq u_\varepsilon \alpha$ . Similarly, we say image  $X$  has an  $M$ -type RC if for all  $\alpha \in \mathcal{P}$ ,  $\ell_M \alpha \leq \varrho_X(\alpha) - \varrho_X(0) \leq u_M \alpha$ . These response curves can be thought of as having a kind of “Lipschitz” condition on their derivatives since the  $\varepsilon$ -type, e.g., are contained within the cone formed by  $\ell_\varepsilon \alpha$  and  $u_\varepsilon \alpha$ . Figure 5 provides examples of response curves that meet the  $\varepsilon/M$  binning criteria. Next, when Alice is forming her bag from this artificial ALASKA II source, she samples (uniformly) an image with  $\varepsilon$ -type RC with probability  $p$  and samples an image with  $M$ -type RC with probability  $1 - p$ . In the previous sections, our binomial model had  $M = 1$  fixed for notational simplicity in the derivations; note that the equations in Section V can be easily generalized to consider arbitrary  $M > \varepsilon$ .

As seen in Figure 6, if we take  $\varepsilon$  (and  $M$ ) as the sample average of the ‘RC slopes at  $\alpha = 0$ ’ of the  $\varepsilon/M$ -type RCs where the slope is estimated using the first three points

$$\hat{b}_X = \frac{1}{2} \left( \frac{\varrho_X(0.05) - \varrho_X(0)}{0.05 - 0} + \frac{\varrho_X(0.1) - \varrho_X(0)}{0.1 - 0} \right), \quad (33)$$

we observe similar behaviors in the size of the bag gain (the maximal drop in detectability), the frequency of local oscillations, and the value of  $B$  where the SRL regime roughly begins (seen by the decay of the amplitude oscillations and increase in detectability for increasing  $B$ ). See the values of ‘avg  $\varepsilon/M$ ’ in Table I for these sample averages of slopes.

Table I  
PARAMETERS FOR NARROW AND WIDE  $\varepsilon/M$  BINNING ON SPLIT2. THE  $\#\varepsilon$  AND  $\#M$  ARE THE NUMBER OF RCS FROM SPLIT2 THAT QUALIFY AS  $\varepsilon/M$ -TYPE. THE AVG  $\varepsilon$  AND AVG  $M$  ARE THE SAMPLE AVERAGES OF THE ESTIMATED SLOPES  $\hat{b}_X$  AT  $\alpha = 0$  (SEE EQ. (33)) FOR  $\varepsilon/M$ -TYPE RCS, RESPECTIVELY.

	$\ell_\varepsilon$	$u_\varepsilon$	$\ell_M$	$u_M$	$\#\varepsilon$	$\#M$	avg $\varepsilon$	avg $M$
Narrow	0	0.08	0.8	3.2	873	1767	0.016	1.564
Wide	0	0.15	0.5	9.5	1358	9149	0.027	3.167

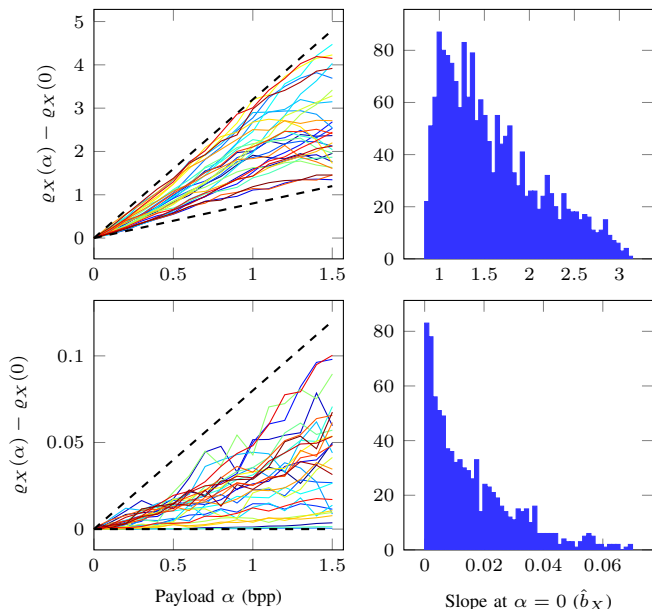


Figure 5. Left column shows examples of response curves  $\varrho_X(\alpha) - \varrho_X(0)$  of bimodal ALASKA II images (narrow binning). Slopes of dashed lines are found in Table I (top row). Right column shows distribution of estimated slopes  $\hat{b}_X$ . Top row is  $M$ -type and bottom row is  $\varepsilon$ -type.

Additionally, Figure 5 shows the distribution of the estimated slopes  $\hat{b}_X$  for both  $\varepsilon/M$ -type RCs.

In Figure 7, observe that there is still a bag loss even when the rejection sampling uses much wider  $\varepsilon/M$  bins. This confirms the robustness of a bag loss occurring even in a source that contains a diverse spectrum of real image response curves (which is very different from binomial linear response curves). If easy-to-steganalyze images are common and hard-to-steganalyze images are rare in an image source, it is important to be aware that a bag loss will likely manifest. Thus, it is important for the steganographer in practice to avoid using bag sizes corresponding to the bag loss to prevent becoming more vulnerable to detection in such cover sources.

## VII. GENERALITY OF THE BAG GAIN

In order for the bag gain to occur, the batch sender must prefer embedding more payload in hard-to-steganalyze images and less payload in easy-to-steganalyze images. In the case of the binomial model, this is equivalent to the batch sender putting more payload in images with near flat response curves. This property holds true for the detector-agnostic IMS, Distortion-Limited Sender (DiLS), and Detectability-Limited Sender (DeLS) studied in [26], as well as the detector-aware

SLS and MDS. The IMS / DiLS / DeLS are not as extreme as the detector-aware senders since they are not designed to explicitly make use of response curves. However, their spreading still correlates with this preference since content-adaptive steganographic schemes put more payload in regions of complex content which give difficulty to detectors. In situations where the steganographers and Warden are knowledge limited as in [30], even a weak preference to embed more in hard-to-steganalyze images (w.r.t. the Warden's detector) can cause the bag gain to manifest.

In this section, we introduce a parametrized family of senders with the parameter controlling how aggressively the sender assigns the payload based on the response curves, including the case when the payload is spread uniformly across all images. By varying this parameter, we can show that the bag gain eventually disappears for sufficiently weak preferences for embedding more payload in harder images.

The Hölder sender can be thought of as a generalization of the MDS (11) as it assigns the following payloads to images:

$$\alpha_i = \frac{rB}{b_i^q \sum_{k=1}^B \frac{1}{b_k^q}}, \quad (34)$$

where  $q \in \mathbb{R}$  is a parameter. For  $q = 2$  and  $q = 1$ , this sender corresponds to the MDS and SLS, respectively. When  $q = 0$ , the payload is spread uniformly across all images.

Following the same steps as in Section V-B, the deflection coefficient for the Hölder sender is

$$\Delta^2(\mathbf{X}) = \begin{cases} \frac{r^2 B^2 \varepsilon^2}{\sigma^2} \left( \frac{C_\varepsilon + \varepsilon^{2q-2}(B-C_\varepsilon)}{(C_\varepsilon + \varepsilon^q(B-C_\varepsilon))^2} \right) & T_q < C_\varepsilon \\ \frac{C_\varepsilon \varepsilon^2 (\log_2 3)^2}{\sigma^2} + \frac{(rB - C_\varepsilon \log_2 3)^2}{\sigma^2 (B - C_\varepsilon)} & T_q \geq C_\varepsilon \end{cases} \quad (35)$$

where  $T_q = \frac{rB}{(1-\varepsilon^q) \log_2 3} - \frac{\varepsilon^q B}{1-\varepsilon^q}$ . Substituting (35) into Eq. (10), we can compute the bag gain  $\gamma$  as given by Eq. (28). Figure 8 shows  $\gamma$  as a function of the exponent  $q$  for a range of the parameters  $p$  (left) and  $\varepsilon$  (right). As  $q$  decreases from  $q = 2$  (MDS), the payload assignment is less polarized and the bag gain starts decreasing. It eventually becomes zero and is always zero for uniform spreading ( $q = 0$ ).

## VIII. RELATIONSHIP TO PRIOR WORK

In this section, we contrast our contribution with previous work [27] that studies optimal bag size in batch steganography and its recent extension to the JPEG domain and pooled steganalysis [1]. We do so in order to highlight the differences and also to briefly discuss possible future directions by combining both approaches. The authors of [27] extended Gaussian Embedding (GE) to batch steganography. Granting the Warden the knowledge of the underlying distributions, a closed-form expression has been derived for the performance of Warden's likelihood ratio test in a specific collection of bags of images. This was used to implement a batch sender with an adaptive batch size called adaBIM.

The first and the main difference between [27] and this paper is the lack of pooled steganalysis. The authors use a performance measure, which is the minimal total detection error  $P_E$  under equal priors of a single-image detector that distinguishes between the cover source and a stego source

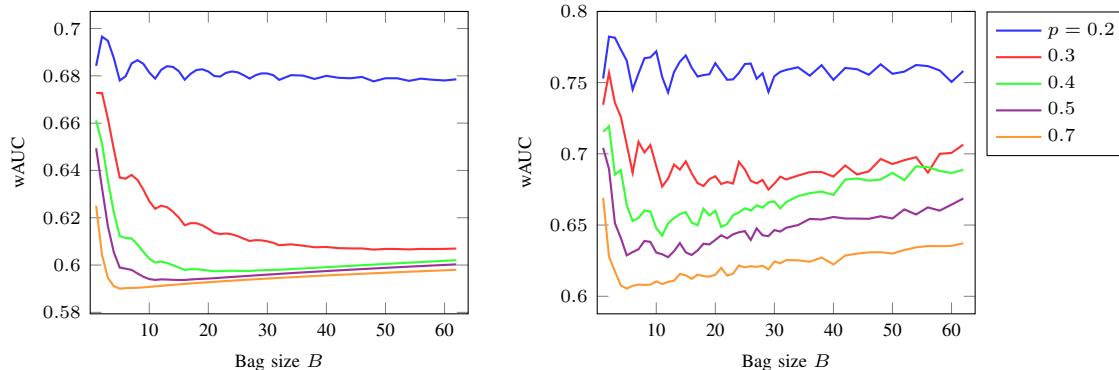


Figure 6. Optimal pooled detector performance for binomial model (left) vs  $\pi_{LRT}$  for bimodal Alaska II (right) using the narrow range parameters given in Table I. The binomial model uses  $\varepsilon = \text{'avg } \varepsilon \text{'}$  and  $M = \text{'avg } M \text{'}$  as given in the table and explained in the text. Rate  $r = 0.3$  bpp is fixed.

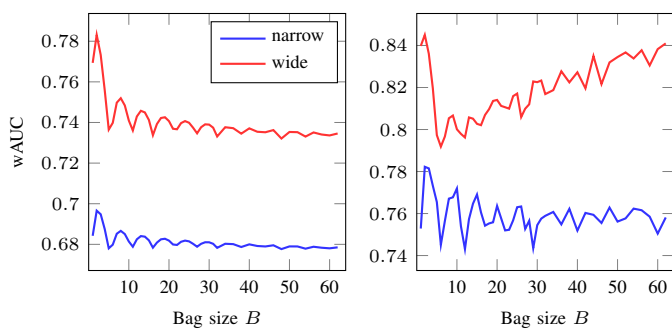


Figure 7. Optimal pooled detector performance for binomial model (left) vs  $\pi_{LRT}$  for bimodal ALASKA II (right) comparing narrow and wider  $\varepsilon/M$  binning. A bag loss still occurs even for wider bins pointing to the robustness of the phenomenon. Parameters  $p = 0.2$  and  $r = 0.3$  are fixed.

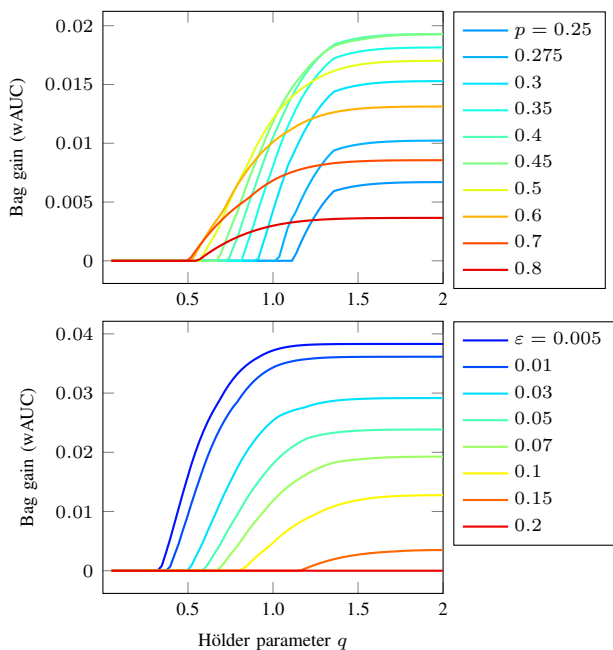


Figure 8. Bag gain measured in wAUC of the bounded capacity model across the family of Hölder senders for  $q \in [0, 2]$  and various  $p$  and  $\varepsilon$ . The left figure has  $\varepsilon = 0.07$  fixed and the right figure has  $p = 0.5$  fixed. Both have  $r = 0.3$  bpp.

whose images contain variable payload “tags” determined by partitioning the dataset into batches and applying GE version of an existing embedding algorithm to the union of all images from the bag to obtain the tags. Furthermore, the effect of bag size in [27] is only studied in asymptotic limits of zero or infinite payloads (Theorem 2). For small payloads, the optimal single-image source detector has highest detection error  $P_E$  when the bag size is equal to the entire image dataset. For large payloads, the highest  $P_E$  occurs when payloads are assigned using bag size 1. This theorem thus only hints at the existence of optimal bag size w.r.t.  $P_E$  and a fixed set of bags. The optimal bag size w.r.t. a single-image detector observed in experiments is merely discussed in words without quantitative results.

This work is extended to JPEG images in [1], where the authors also derive the optimal pooled detector for a given bag. To assess the performance of the pooled detector across bags, however, one needs to consider the variability of bags, which would necessitate adopting a meta-model on the source. In the case of the GE as studied in [27], [1], it would likely have to be the distribution of the product of cover image pixel variances, which opens the possibility to use, e.g., a similar binomial model within the context of GE. We plan to investigate this direction in the future.

In contrast, the approach taken in this paper allowed us to relate all essential aspects of a steganographic channel—the cover source diversity, detector response, payload, and bag size—to security under pooled steganalysis. We also believe that working with detector output models leads to a tighter correspondence between the detectability derived from the model and the one obtained experimentally. After all, the model correctly predicts completely new phenomena, such as the bag loss and local oscillations in the small bag regime.

## IX. CONCLUSIONS

In batch steganography, the secret payload is spread among multiple cover images forming a bag. Within the context of content-adaptive steganography, many batch senders were proposed and studied in the past, such as the image merging sender [26], [27] and the deflection/distortion limited senders [26], as well as two detector-aware senders, the shift limited sender and the minimum deflection sender [30]. When

a fixed relative payload is communicated in each bag, batch senders that embed larger payloads in difficult-to-steganalyze images and smaller payloads in easy images exhibit similar trends in terms of detectability vs. the bag size. In this paper, we analyze these trends from the simplest model that captures their essence by considering only two types of images that are “easy” and “difficult” to steganalyze. While the trends depend on the cover source diversity, detector response characteristics, batch sender, and the communication rate, our work offers a simple intuitive explanation.

The square root law states that, when the sender maintains a fixed communication rate, the detectability increases with the number of images sent (the bag size). Due to variability of images across bags, however, the law engages differently in bags of different sizes. Assuming that difficult images that can hold large payloads are rare, as the bag size increases, initially the detectability as measured with pooled detectors *increases* due to square root law because only a small fraction of bags contains the difficult images that can carry large payloads without triggering a detector – the square root law thus engages based on embedding primarily in easy images. Once the bag size becomes large enough to contain difficult images with high probability, they hold most of the payload and the detectability begins to *decrease*. Due to the square root law, the detectability eventually levels off, reaching a global minimum, and once more increases but at a speed *slower* than the initial rise depending on the ratio of easy and difficult images in the cover source and the communication rate. The maximum initial rise in detectability is called the *bag loss* while the global minimum corresponds to a *bag gain*. Both phenomena essentially manifest because the average statistical make up of bags differs between small and large bags, which affects how the square root law engages.

While the bag gain was observed experimentally in previous art [30], it was a mere experimental fact that was left unexplained. Our work provides theoretical insight into the manifestation of the bag gain and quantifies how it depends on cover source diversity, detector response, batch sender, and communication rate. The predicted trends closely match experiments with real images. The predicted bag loss, together with some higher-order oscillations, are experimentally confirmed in datasets with suitable diversity. Furthermore, we provide evidence that these phenomena manifest for batch senders that generally assign payloads based on detectability of embedding in individual images sufficiently strongly as bag loss and gain are not observed for uniform batch senders.

On the practical side, our work shows that it is important to be aware of the existence of the bag gain and bag loss for practitioners who need to avoid combinations of bag sizes and communication rates that lead to bag loss and select the bag size that corresponds to bag gain as this will make the covert communication less detectable to an adversary.

## APPENDIX

### A. ROC for pooled detector

A pooled detector makes a decision on bags—either it contains cover or stego images. Since the images from each

bag are randomly selected from a cover source, some bags will be easier to detect than others, depending on the value of the deflection coefficient  $\Delta^2$ . In this section, we derive an expression for the ROC of the pooled detector over bags based on the distribution of the deflection coefficient.

For a fixed false-alarm  $P_{\text{FA}}$ , the probability of correct stego bag detection is

$$P_{\text{D}}(P_{\text{FA}}) = \mathbb{E}[Q(Q^{-1}(P_{\text{FA}}) - \Delta)], \quad (36)$$

the expectation taken over bags. In this paper,  $\Delta$  is discrete, attaining values from a finite set  $\mathcal{D}$ . The derivation below, however, is also valid for a continuous-valued  $\Delta$ . Let  $p_{\Delta}(x)$ ,  $x \in \mathcal{D}$ , be the probability mass function of  $\Delta$  and let  $\mu = \mathbb{E}[\Delta]$ . Then, using Taylor expansion of  $Q(x)$  at  $Q^{-1}(P_{\text{FA}}) - \mu$  with the Lagrange form for the remainder, the expected ROC (36) can be written as

$$\begin{aligned} P_{\text{D}}(P_{\text{FA}}) &= \sum_{x \in \mathcal{D}} Q(Q^{-1}(P_{\text{FA}}) - x) p_{\Delta}(x) \\ &= \sum_{x \in \mathcal{D}} \left[ \sum_{k=0}^{n-1} \frac{(\mu - x)^k}{k!} Q^{(k)}(Q^{-1}(P_{\text{FA}}) - \mu) \right. \\ &\quad \left. + \frac{(\mu - x)^n}{n!} Q^{(n)}(Q^{-1}(P_{\text{FA}}) - x^*) \right] p_{\Delta}(x) \\ &= \sum_{k=0}^{n-1} \frac{(-1)^k c_k}{k!} Q^{(k)}(Q^{-1}(P_{\text{FA}}) - \mu) + \quad (37) \\ &\quad \underbrace{\sum_{x \in \mathcal{D}} \frac{(\mu - x)^n}{n!} Q^{(n)}(Q^{-1}(P_{\text{FA}}) - x^*) p_{\Delta}(x)}_{R_n}, \quad (38) \end{aligned}$$

where  $c_k$  is the  $k$ th central moment of  $\Delta$  and  $x^* \in (\mu, x)$ . Next, we leverage the following bound on the  $n$ th derivative of the  $Q$  function

$$|Q^{(n)}(x)| \leq \sqrt{\frac{n!}{2\pi}} \text{ for all } x, \quad (39)$$

which follows from the fact that  $Q^{(n)}(x) = \frac{1}{\sqrt{2\pi}} p_n(x) e^{-x^2/2}$ , where  $p_n(x)$  is the statistician’s Hermite polynomial (the physicist’s Hermite polynomial is  $H_n(x) = 2^{n/2} p_n(\sqrt{2}x)$ ), and Cramér inequality [14] for Hermite function defined using physicist’s Hermite polynomials  $\Psi_n(x) = (2^n n! \sqrt{\pi})^{-1/2} e^{-x^2/2} H_n(x) \leq \pi^{-1/4}$  for all  $x$  and all  $n$ . Hence

$$|R_n| \leq \frac{c_n}{n!} \sqrt{\frac{n!}{2\pi}} = \frac{c_n}{\sqrt{2\pi n!}}. \quad (40)$$

### B. General form of the MDS

Let  $r \in [0, \log_2 3]$  be a chosen embedding rate in bpp, and assume the response curves are linear with slopes  $b_i \in [0, \infty)$  for all  $i$ . Optimal payloads for the MDS are found by minimizing  $\Delta^2(\mathbf{X})$  under constraints  $\sum_{i=1}^B \alpha_i = rB$  and  $\alpha_i \in [0, A_i] \forall i$  where  $A_i \leq \log_2 3$  is the embedding

capacity of the  $i$ th image (accounting for wet pixels [9]). The Lagrangian has the form

$$\begin{aligned} \mathcal{L} = & \sum_{i=1}^B b_i^2 \alpha_i^2 - \lambda \left( \sum_{i=1}^B \alpha_i - rB \right) \\ & - \sum_{i=1}^B \ell_i \alpha_i - \sum_{i=1}^B u_i (\alpha_i - A_i), \end{aligned} \quad (41)$$

where  $\ell_i$  and  $u_i$  are KKT multipliers that satisfy the lower and upper inequality constraints on  $\alpha_i$ , respectively. To be a stationary point, the tuple  $(\alpha_1, \dots, \alpha_B)$  must satisfy

$$\alpha_i = 0, \frac{\lambda}{2b_i^2}, \text{ or } A_i; \quad \forall i. \quad (42)$$

Let  $\mathcal{L}$  and  $\mathcal{U}$  denote the sets of indices for which  $\alpha_i = 0$  or  $A_i$ , respectively. Let  $\mathcal{I} = (\mathcal{L} \cup \mathcal{U})^c$  be the set of remaining indices where  $0 < \alpha_i < A_i$ . From the payload constraint

$$\begin{aligned} rB &= \sum_{k \in \mathcal{L}} 0 + \sum_{k \in \mathcal{I}} \frac{\lambda}{2b_k^2} + \sum_{k \in \mathcal{U}} A_k \\ \Rightarrow \lambda &= \frac{rB - \sum_{k \in \mathcal{U}} A_k}{\frac{1}{2} \sum_{k \in \mathcal{I}} \frac{1}{b_k^2}} \\ \Rightarrow \alpha_i &= \frac{rB - \sum_{k \in \mathcal{U}} A_k}{b_i^2 \sum_{k \in \mathcal{I}} \frac{1}{b_k^2}}, \end{aligned} \quad (43)$$

for all  $i \in \mathcal{I}$ . The optimal payload is found numerically by searching over the combinations of  $\mathcal{L}$ ,  $\mathcal{I}$ , and  $\mathcal{U}$ .

Note that when  $A_i = \infty$  for all  $i$  (unbounded embedding capacity), we have  $\mathcal{U} = \mathcal{L} = \emptyset$  and (43) simplifies to

$$\alpha_i = \frac{rB}{b_i^2 \sum_{k=1}^B \frac{1}{b_k^2}}. \quad (44)$$

## REFERENCES

- [1] M. Aloraini, M. Sharifzadeh, and D. Schonfeld. Quantized gaussian JPEG steganography and pool steganalysis. *IEEE Access*, 10:38031–38044, 2022.
- [2] P. Bas. Steganography via cover-source switching. In *2016 IEEE International Workshop on Information Forensics and Security (WIFS)*, pages 1–6, December 4–7 2016.
- [3] P. Bas and T. Furon. BOWS-2. <http://bows2.ec-lille.fr>, July 2007.
- [4] M. Boroumand, M. Chen, and J. Fridrich. Deep residual network for steganalysis of digital images. *IEEE Transactions on Information Forensics and Security*, 14(5):1181–1193, May 2019.
- [5] J. Butora, Y. Yousfi, and J. Fridrich. How to pretrain for steganalysis. In D. Borghys and P. Bas, editors, *The 9th ACM Workshop on Information Hiding and Multimedia Security*, Brussels, Belgium, 2021. ACM Press.
- [6] R. Cogranne, Q. Giboulot, and P. Bas. ALASKA-2: Challenging academic research on steganalysis with realistic images. In *IEEE International Workshop on Information Forensics and Security*, New York, NY, December 6–11, 2020.
- [7] T. Filler, T. Pevný, and P. Bas. BOSS (Break Our Steganography System). <http://www.agents.cz/boss>, July 2010.
- [8] J. Fridrich. *Steganography in Digital Media: Principles, Algorithms, and Applications*. Cambridge University Press, 2009.
- [9] J. Fridrich, M. Goljan, D. Soukal, and P. Lisoněk. Writing on wet paper. In T. Kalker and P. Moulin, editors, *IEEE Transactions on Signal Processing, Special Issue on Media Security*, volume 53, pages 3923–3935, October 2005. (journal version).
- [10] Q. Giboulot, P. Bas, and R. Cogranne. Multivariate side-informed Gaussian embedding minimizing statistical detectability. *IEEE Transactions on Information Forensics and Security*, 17:1841–1854, 2022.
- [11] Q. Giboulot, R. Cogranne, and P. Bas. Detectability-based JPEG steganography modeling the processing pipeline: The noise-content trade-off. *IEEE Transactions on Information Forensics and Security*, 16:2202–2217, 2021.
- [12] V. Holub, J. Fridrich, and T. Denemark. Universal distortion design for steganography in an arbitrary domain. *EURASIP Journal on Information Security, Special Issue on Revised Selected Papers of the 1st ACM IH and MMS Workshop*, 2014:1, 2014.
- [13] X. Hu, J. Ni, W. Zhang, and J. Huang. Efficient JPEG batch steganography using intrinsic energy of image contents. *IEEE Transactions on Information Forensics and Security*, 16:4544–4558, 2021.
- [14] J. Indritz. An inequality for Hermite polynomials. 12(6):981–983, 1961.
- [15] J. R. Janesick. *Scientific Charge-Coupled Devices*, volume Monograph PM83. Washington, DC: SPIE Press - The International Society for Optical Engineering, January 2001.
- [16] A. D. Ker. Batch steganography and pooled steganalysis. In J. L. Camenisch, C. S. Collberg, N. F. Johnson, and P. Sallee, editors, *Information Hiding, 8th International Workshop*, volume 4437 of Lecture Notes in Computer Science, pages 265–281, Alexandria, VA, July 10–12, 2006. Springer-Verlag, New York.
- [17] A. D. Ker. Batch steganography and the threshold game. In E. J. Delp and P. W. Wong, editors, *Proceedings SPIE, Electronic Imaging, Security, Steganography, and Watermarking of Multimedia Contents IX*, volume 6505, pages 04 1–13, San Jose, CA, January 29–February 1, 2007.
- [18] A. D. Ker. A capacity result for batch steganography. *IEEE Signal Processing Letters*, 14(8):525–528, 2007.
- [19] A. D. Ker. Perturbation hiding and the batch steganography problem. In K. Solanki, K. Sullivan, and U. Madhow, editors, *Information Hiding, 10th International Workshop*, volume 5284 of Lecture Notes in Computer Science, pages 45–59, Santa Barbara, CA, June 19–21, 2008. Springer-Verlag, New York.
- [20] A. D. Ker. Locally square distortion and batch steganographic capacity. *International Journal of Digital Crime and Forensics*, 1(1):29–44, 2009.
- [21] A. D. Ker. The square root law of steganography. In M. Stamm, M. Kirchner, and S. Voloshynovskiy, editors, *The 5th ACM Workshop on Information Hiding and Multimedia Security*, Philadelphia, PA, June 20–22, 2017. ACM Press.
- [22] A. D. Ker and Tomas Pevný. Batch steganography in the real world. In J. Dittmann, S. Craver, and S. Katzenbeisser, editors, *Proceedings of the 14th ACM Multimedia & Security Workshop*, pages 1–10, Coventry, UK, September 6–7, 2012.
- [23] B. Li, M. Wang, and J. Huang. A new cost function for spatial image steganography. In *Proceedings IEEE, International Conference on Image Processing, ICIP*, Paris, France, October 27–30, 2014.
- [24] L. Li, W. Zhang, C. Qin, K. Chen, W. Zhou, and N. Yu. Adversarial batch image steganography against CNN-based pooled steganalysis. *Signal Processing*, 181:107920–107920, 2021.
- [25] T. Pevný and I. Nikolaev. Optimizing pooling function for pooled steganalysis. In *IEEE International Workshop on Information Forensics and Security*, pages 1–6, Rome, Italy, November 16–19, 2015.
- [26] V. Sedighi, R. Cogranne, and J. Fridrich. Practical strategies for content-adaptive batch steganography and pooled steganalysis. In *Proceedings IEEE, International Conference on Acoustics, Speech, and Signal Processing*, March 5–9, 2017.
- [27] M. Sharifzadeh, M. Aloraini, and D. Schonfeld. Adaptive batch size image merging steganography and quantized Gaussian image steganography. *IEEE Transactions on Information Forensics and Security*, 15:867–879, 2020.
- [28] T. Thai, R. Cogranne, and F. Retraint. Statistical model of quantized DCT coefficients: Application in the steganalysis of Jsteg algorithm. *Image Processing, IEEE Transactions on*, 23(5):1–14, May 2014.
- [29] Thanh Hai Thai, R. Cogranne, and F. Retraint. Camera model identification based on the heteroscedastic noise model. *Image Processing, IEEE Transactions on*, 23(1):250–263, Jan 2014.
- [30] Y. Yousfi, E. Dworetzky, and J. Fridrich. Detector-informed batch steganography and pooled steganalysis. In J. Butora, B. Tondi, and C. Veilhauer, editors, *The 10th ACM Workshop on Information Hiding and Multimedia Security*, Santa Barbara, CA, 2022. ACM Press.
- [31] A. Zakaria, M. Chaumont, and G. Subsol. Pooled steganalysis in JPEG: how to deal with the spreading strategy? In *IEEE International Workshop on Information Forensics and Security (WIFS)*, pages 1–6, 2019.
- [32] N. Zhong, Z. Qian, Z. Wang, X. Zhang, and X. Li. Batch steganography via generative network. *IEEE Transactions on Circuits and Systems for Video Technology*, 31(1):88–97, January 2021.