# DESIGNING STEGANOGRAPHIC DISTORTION USING DIRECTIONAL FILTERS

*Vojtěch Holub and Jessica Fridrich*

Department of ECE, SUNY Binghamton, NY, USA
{vholub1, fridrich}@binghamton.edu

## ABSTRACT

This paper presents a new approach to defining additive steganographic distortion in the spatial domain. The change in the output of directional high-pass filters after changing one pixel is weighted and then aggregated using the reciprocal Hölder norm to define the individual pixel costs. In contrast to other adaptive embedding schemes, the aggregation rule is designed to force the embedding changes to highly textured or noisy regions and to avoid clean edges. Consequently, the new embedding scheme appears markedly more resistant to steganalysis using rich models. The actual embedding algorithm is realized using syndrome-trellis codes to minimize the expected distortion for a given payload.

## 1. INTRODUCTION

Designing steganographic algorithms for empirical cover sources, such as digital images, is very challenging due to the fundamental lack of accurate models. The most successful approach today avoids estimating the cover source distribution because this task is infeasible for complex and highly non stationary sources. Instead, the steganography problem is formulated as source coding with fidelity constraint [2] – the sender embeds her message while minimizing an appropriately defined distortion. Practical algorithms that embed near the theoretical payload–distortion bound are available for a very general class of distortion functions [4, 2]. Within this framework, the only task left to the sender is essentially the design of the distortion function.

In an attempt to relate distortion with statistical detectability, the authors of [3] parametrized the distortion function and then searched for such values of the parameters that gave the smallest detectability evaluated as a margin between classes within a selected feature space (cover model). However, unless the cover model is a complete statistical descriptor of the empirical source [10], such optimized schemes may, paradoxically, end up being more detectable if the Warden designs the detector "outside of the model" [11], which brings us back to the main and rather difficult problem – modeling the source.

All of today's most secure steganographic schemes for digital images use heuristically defined distortion functions that constrain the embedding changes to those parts of the image that are difficult to model (e.g., complex textures or "noisy" areas). In the JPEG domain, by far the most successful approach is built around distortion functions that measure distortion w.r.t. the raw, uncompressed image [9, 15, 16]. A natural way to define the distortion function in the spatial domain is to assign pixel costs by measuring the impact of changing each pixel in a feature (model) space using a weighted norm. Making the weights dependent on the pixel's local neighborhood introduces desirable content adaptivity. An example of this approach is the embedding algorithm HUGO [14], which employs the SPAM feature model. To the best knowledge of the authors, and based on the recent steganalysis study [6], HUGO is currently the most secure algorithm for embedding in the spatial domain even though its secure payload has been substantially lowered by modern attacks initiated during the BOSS competition [5] that employ high-dimensional rich models.

In this paper, we approach the task of building distortion functions in the spatial domain using a different strategy. Instead of using a weighted norm in some steganalytic model to compute the pixel costs, we employ a bank of directional high-pass filters to obtain the so-called directional residuals, which are related to the predictability of the pixel in a certain direction. By measuring the impact of embedding on every directional residual and by suitably aggregating these impacts, we force the embedding cost to be high where the content is predictable in at least one direction (smooth areas and along edges) and low where the content is unpredictable in every direction (e.g., in textured or noisy areas). The resulting algorithm thus becomes highly adaptive and better resists steganalysis using rich models.

After introducing basic notation in Section 2, we list three steganographic methods with which new schemes will be compared using an empirical measure of security. In Section 3, we describe the distortion function, including the filter banks for computing the directional residuals and the

aggregation rule. The purpose of the exploratory analysis of Section 4 is to assess the effect of various design elements on security and select the setting that provides the highest empirical security. In Section 5, we subject the new scheme to steganalysis in the wavelet domain where the embedding costs are computed. The paper is concluded in Section 6.

## 2. PRELIMINARIES

Capital and lower-case boldface symbols stand for matrices and vectors, respectively. The symbols $\mathbf{X} = (X_{ij}), \mathbf{Y} = (Y_{ij}) \in \{0, \ldots, 255\}^{n_1 \times n_2}$ will always be used an 8-bit gray-scale cover (and the corresponding stego) image with $n_1 \times n_2$ pixels. For matrix $\mathbf{X}$, $\mathbf{X}^\mathrm{T}$ is its transpose, $\mathbf{X}^\frown$ is $\mathbf{X}$ rotated by 180 degrees, and $|\mathbf{X}|$ is the matrix of absolute values.

### 2.1. Empirical security

All experiments are conducted on BOSSbase ver. 1.0 [5] with 10 000 images. The steganographic security is evaluated empirically using binary classifiers trained on a given cover source and its stego version embedded with a fixed payload. With the exception of Section 5, we use the Spatial Rich Model (SRM) [6] consisting of 106 symmetrized sub-models with a total dimension of $34,671$. All classifiers were implemented using the ensemble [12] with Fisher linear discriminants as base learners. Security is quantified using the ensemble's "out-of-bag" (OOB) error $E_{\mathrm{OOB}}$, which is an unbiased estimate of the testing error "averaged" over multiple bootstrap samples of the image source during training [12].

### 2.2. Steganography methods

We compare the proposed methods with HUGO, the Edge Adaptive (EA) algorithm [13], and Least Significant Bit Matching (LSBM). We used the embedding simulator [5] for HUGO operating at the theoretical payload–distortion bound with default settings $\gamma = 1$, $\sigma = 1$, and the switch --T with $T = 255$ to remove the weakness reported in [11]. LSBM was simulated at the ternary entropy bound. The code for the EA algorithm with its custom coding scheme was obtained from the authors.

## 3. DISTORTION FUNCTION DESIGN

We restrict our design to additive distortion in the form:

$$D(\mathbf{X}, \mathbf{Y}) = \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} \rho_{ij}(\mathbf{X}, Y_{ij})|X_{ij} - Y_{ij}|, \quad (1)$$

where $\rho_{ij}$ are the costs of changing pixel $X_{ij}$ to $Y_{ij}$. The additivity means that we do not consider the effects of individual embedding changes influencing each other. We opted



**Table 1**. Filter banks used in this paper.

for this mainly to simplify the design. Since the embedding algorithm will be forced to concentrate the embedding modifications into highly textured/noisy areas, using the Gibbs construction [2] with non-additive distortion functions may have an additional beneficial impact on security. The authors contemplate investigating this direction as part of future research.

Having defined the pixel costs $\rho_{ij}$, embedding a (pseudo) random sequence of bits with minimal expected distortion (1) is equivalent to source coding with a fidelity criterion. A practical algorithm, based on Syndrome-Trellis Codes (STCs), that embeds near the payload–distortion bound was proposed in [4]. It works in the dual domain to better cover the range of small payloads typically needed for steganography. The STC Toolbox, which we also use in this paper to implement all our schemes, can be downloaded from http://dde.binghamton.edu/download/syndrome/.

### 3.1. Directional Filters

As already reported by its authors, the distortion function of HUGO concentrates the embedding changes primarily in textures and edges. However, the content along an edge can usually be well modeled using locally polynomial models, which aids the detection [7, 6, 8]. Thus, whenever possible the embedding algorithm should embed into textured/noisy areas that are not easily modellable in any direction. To this end, we evaluate the smoothness in multiple directions using a filter bank $\mathcal{B}_n = \{\mathbf{K}^{(1)}, \ldots, \mathbf{K}^{(n)}\}$ consisting of $n$ multiple directional high-pass filters represented by their kernels normalized so that all $L_2$-norms $\left\|\mathbf{K}^{(k)}\right\|_2$ are the same. The $k$-th residual $\mathbf{R}^{(k)}$, $k = 1, \ldots, n$, is computed as $\mathbf{R}^{(k)} = \mathbf{K}^{(k)} \star \mathbf{X}$, where '$\star$' is a convolution mirror-padded so that $\mathbf{R}^{(k)}$ has

again $n_1 \times n_2$ elements. (The mirror-padding prevents introducing embedding artifacts at image boundary.) If the residual values $\mathbf{R}_{ij}^{(k)}$ are large for some $ij$ and for all $k$, it means that the local content at pixel $x_{ij}$ is not smooth in any direction and thus difficult to model.

Since we want to detect edges in all directions, it is natural to use established edge detectors for the filter banks (see Table 1). The non-directional 'KB' filter [1] is often used in steganalysis, while the Sobel operator is a common edge detector. Wavelet-based Directional Filter Banks 'WDFB-H' and 'WDFB-D' use the Haar and Daubechies 8-tap wavelets. The computation of the residual coincides with the first-level wavelet decomposition with no decimation. The wavelet banks consist of three filters, $\mathbf{K}^{(1)}, \mathbf{K}^{(2)}, \mathbf{K}^{(3)}$, using which the LH, HL, and HH directional residuals are obtained. Given the wavelet's 1-D low-pass decomposition filter $\mathbf{h}$ and a high-pass decomposition filter $\mathbf{g}$, the 2-D directional filters are computed as shown in Table 1.

### 3.2. Aggregating embedding suitability

The embedding should prefer changing large values of directional residuals, where the textures and edges are, and preserve the small values, where the content is predictable. One way to achieve this is to weigh the difference between $\mathbf{R}^{(k)}$ and the same residual after changing only one pixel at $ij$ (denoted $\mathbf{R}_{[ij]}^{(k)}$) by the wavelet coefficient itself:

$$\xi_{ij}^{(k)} = \left| \mathbf{R}^{(k)} \right| \star \left| \mathbf{R}^{(k)} - \mathbf{R}_{[ij]}^{(k)} \right|^{\curvearrowright} \overset{(a)}{=} \left| \mathbf{R}^{(k)} \right| \star \left| \mathbf{K}^{(k)} \right|^{\curvearrowright}. \quad (2)$$

The quantity $\xi_{ij}^{(k)}$, which we call embedding "suitability," is formally a correlation between the absolute value of the cover residual with the absolute value of the residual change. Since $\mathbf{R}^{(k)} - \mathbf{R}_{[ij]}^{(k)}$ is the spatially shifted directional filter $\mathbf{K}^{(k)}$, $\xi_{ij}^{(k)}$ can be computed for all pixels at once (equality $(a)$).

Next, we compute the embedding costs $\rho_{ij}$ by aggregating all suitabilities $\xi_{ij}^{(k)}$, $k = 1, \dots, n$. Since we wish to restrict the embedding changes to those pixels with complex content in every direction, the aggregation rule $\rho : \mathbb{R}^n \longrightarrow \mathbb{R}_0^+$, $\rho_{ij} = \rho(\xi_{ij}^{(1)}, \dots, \xi_{ij}^{(n)})$ is required to have the following properties:

A1. The larger the values of $|\xi_{ij}^{(k)}|$, the smaller the $\rho_{ij}$ should be.

A2. If there exists $k \in \{1, \dots, n\}$ such that $\xi_{ij}^{(k)} = 0$, then $\rho_{ij} = +\infty$.

A simple function that meets both requirements is the reciprocal Hölder norm with $p < 0$:

$$\rho_{ij}^{(p)} = \left( \sum_{k=1}^{n} |\xi_{ij}^{(k)}|^p \right)^{-\frac{1}{p}}. \quad (3)$$



**Fig. 1**. $E_{\text{OOB}}$ as a function of the Hölder-norm parameter $p$ when embedding at 0.4 bpp with the WDFB.



**Fig. 2**. Evaluating the security of several different filter banks using the OOB estimate of the testing error, $E_{\text{OOB}}$.

We restrict the embedding changes to $\pm 1$, $|X_{ij} - Y_{ij}| = 1$. Note that due to the absolute value in (2), both changes result in the same embedding cost, which allows us to use the more powerful multi-layered version of STCs [4] also available in the STC Toolbox (see Section 4.4 for a discussion of how different coding schemes affect the security).

### 4. EXPERIMENTS

In this section, we first assess how various design parameters, such as $p$, the filter bank, and coding, affect security. Then, the most secure setting is identified and compared with HUGO, EA, and LSBM.

### 4.1. Aggregation rule

To obtain an insight as to which value of $p$ should be used in the aggregation rule (3), in Fig. 1 we plot $E_{\text{OOB}}(p)$ for the WDFB when embedding the payload of 0.4 bpp (bits per pixel). While for $p < 0$ the security appears almost constant, for $p > 0$ the requirement A2 is longer valid – the costs at

smooth edges decrease, which lowers the security. A similar dependence of security on $p$ was observed for other filter banks. Thus, for concreteness and simplicity, we fix the value of the parameter $p$ to $p = -1$.

### 4.2. Assessing filter banks

Fig. 2 shows the OOB error estimate for filter banks listed in Table 1. Among them, the WDFB-D achieves the best steganographic security. We call this embedding algorithm WOW (Wavelet Obtained Weights). All the other filters achieve comparable security with HUGO, even the WDFB-H with support of size only $2 \times 2$.

Encouraged with the success of the Daubechies 8-tap wavelet-based filter bank, we experimented with several other wavelet bases, including the Biorthogonal 44 wavelets, only to achieve very similar results in terms of the $E_{\mathrm{OOB}}$.

### 4.3. Comparison to prior art

Fig. 3 shows the comparison between WOW and three other algorithms using the SRM model (left) and a model constructed using dependencies in the wavelet domain (see Section 5 for more details). The improvement over HUGO is especially apparent for large payloads – at 0.5 bpp, the $E_{\mathrm{OOB}}$ of WOW is almost twice as high as that of HUGO.

### 4.4. The effect of coding

As already mentioned in Section 3.2, since the costs (3) do not depend on the direction of the embedding change, WOW can use the ternary multi-layered version of STCs. Fig. 4 shows that the gain of using the ternary STCs over their binary version is quite significant. At the same time, the coding loss of STCs w.r.t. optimal embedding operating at the payload–distortion bound is rather small.

The last comment above might suggest that HUGO might be improved using ternary embedding instead of binary. However, since HUGO embeds only in the direction of smaller distortion and allows interaction among modifications, it is not clear how to implement ternary embedding and what the security impact would be.

### 4.5. WOW adaptivity

In Fig. 5, we contrast the placement of embedding changes for HUGO and for WOW. The selected cover image has numerous horizontal and vertical edges and also some textured areas. While HUGO embeds with high probability into the pillar edges as well as the horizontal lines above the pillars, WOW embeds solely into the textured areas as dictated by the aggregation rule (3).



**Fig. 4**. OOB error estimate of the testing error, $E_{\mathrm{OOB}}$, for WOW implemented using binary and ternary STCs versus simulated optimal embedding. Note the large gain of ternary STCs versus their binary version. Also note that the coding loss is quite small.



**Fig. 5**. Embedding probability for payload $0.4$ bpp using HUGO (bottom left) and WOW (bottom right) for a $128 \times 128$ grayscale cover image (top).

**Fig. 3**. Comparing statistical detectability of WOW and three state-of-the-art embedding algorithms using the SRM (left) and wavelet-domain dependencies (right).

## 5. STEGANALYSIS IN WAVELET DOMAIN

The most successful steganalysis attacks have always been built in the embedding domain. Although WOW embeds in the spatial domain, which is well covered by the SRM, the costs are computed in a transform domain. The goal of this section is to investigate whether WOW can be attacked in the wavelet domain by forming features that capture dependencies among wavelet coefficients.

Inspired by how steganalysis features are built in the JPEG domain, we explored the following four logical possibilities graphically shown in Fig. 6: 4-D co-occurrence matrices built from four consecutive wavelet coefficients exploiting dependencies a) intra band (IaB), b) inter level, c) a mix of intra level and inter level, and d) intra level, inter band. The IaB features are similar in spirit to the SRM provided the wavelet coefficients are interpreted as noise residuals. Since the IaB features were much more successful in detecting WOW when compared to the other three possibilities b)–d), we only provide detailed discussion for case a).

The coefficients were computed using the standard undecimated discrete wavelet decomposition with the Daubechies 8 wavelet (to steganalyze in the domain where WOW computes its costs).

Let $\mathbf{S}^{(l,s)} = \{c_{ij}^{(l,s)}\}$, $l \in \{1,2,3\}$, $s \in \{\text{LH}, \text{HL}, \text{HH}\}$, be the undecimated $s$th subband in the $l$th level of the wavelet transform. Assuming that $n_1 = 2^{k_1}$, $n_2 = 2^{k_2}$ for some $k_1, k_2 \in \mathbb{Z}$, the range of subscripts for $c_{ij}^{(l,s)}$ is $i = \{1, \ldots, 2^{k_1-l+1}\}$ and $j = \{1, \ldots, 2^{k_2-l+1}\}$.

The steganalytic features are four-dimensional co-occurrence matrices formed by groups of four horizontally and vertically adjacent coefficients after truncation and quantization to a

finite dynamic range, $c_{ij}^{(l,s)} \leftarrow \text{round}\left(\text{trunc}_T(c_{ij}^{(l,s)}/q)\right)$, where $q$ is a quantization step and $\text{trunc}_T(x) = x$ for $x \in [-T, T]$, $\text{trunc}_T(x) = T \cdot \text{sign}(x)$ otherwise. The horizontal co-occurrence matrix is denoted as $\mathbf{C}_{\mathbf{d}}^{(\text{h},l,s)}$, $\mathbf{d} = (d_1, \ldots, d_4) \in \{-T, \ldots, T\}^4$, $\mathbf{C}_{\mathbf{d}}^{(\text{h},l,s)} = \left\{(i,j) \Big| c_{ij}^{(l,s)} = d_1, c_{i,j+1}^{(l,s)} = d_2, c_{i,j+2}^{(l,s)} = d_3, c_{i,j+3}^{(l,s)} = d_4\right\}$, with the vertical matrix $\mathbf{C}_{\mathbf{d}}^{(\text{v},l,s)}$ defined analogically.

We built three co-occurrence matrices for each level $l \in \{1,2,3\}$. In Fig. 6a), denoted by a triangle is the co-occurrence $\mathbf{C}_{\mathbf{d}}^{(1,l)} \triangleq \mathbf{C}_{\mathbf{d}}^{(\text{h},l,\text{LH})} + \mathbf{C}_{\mathbf{d}}^{(\text{v},l,\text{HL})}$. The squares correspond to $\mathbf{C}_{\mathbf{d}}^{(2,l)} \triangleq \mathbf{C}_{\mathbf{d}}^{(\text{v},l,\text{LH})} + \mathbf{C}_{\mathbf{d}}^{(\text{h},l,\text{HL})}$, while the circles mark $\mathbf{C}_{\mathbf{d}}^{(3,l)} \triangleq \mathbf{C}_{\mathbf{d}}^{(\text{h},l,\text{HH})} + \mathbf{C}_{\mathbf{d}}^{(\text{v},l,\text{HH})}$.

In this paper, we used $T = 2$, which gave each co-occurrence matrix $\mathbf{C}_{\mathbf{d}}^{(i,l)}$ the dimensionality of $(2T+1)^4 = 625$. Since $i \in \{1,2,3\}$ and $l \in \{1,2,3\}$, the total number of co-occurrence matrices is 9, giving the final feature vector a dimensionality of $625 \times 9 = 5,625$. A brief study on the effect of the quantization step $q \in [0.2, 5]$ on $E_{\text{OOB}}$ showed that the best performance was usually obtained for $q \approx 1$. Thus, in all our experiments, we set $q = 1$.

Fig. 3 (right) shows the results of steganalysis using the IaB features. WOW still achieves better security than any other tested method. The overall detection performance of the IaB features is, however, inferior to the SRM (left).

## 6. CONCLUSION

This paper confirms what has been suspected before – restricting the embedding changes to textures while avoiding "clean"

**Fig. 6.** Four types of groups of wavelet coefficients from which 4-D co-occurrence matrices were built to steganalyze WOW.

edges greatly improves steganographic security. This high level of adaptivity was achieved through a novel design of the steganographic distortion function. First, a directional filter bank is used to detect edges in local neighborhoods of each pixel. Then the changes in the residuals caused by embedding are weighted and aggregated using a special rule designed to output a low embedding cost only when the local content is not smooth in any direction.

According to our experiments, 2-D wavelet decomposition filters provide the highest level of steganographic security measured empirically for a given image source (database) and classifiers operating in high-dimensional feature spaces (rich image models). The proposed algorithm, WOW, outperforms the current state-of-the-art HUGO by a significant margin especially for large payloads.

Further potential improvement is possible by employing better directional filter banks and by using non-additive distortion to model interaction of embedding changes.

## 7. REFERENCES

[1] R. Böhme. *Improved Statistical Steganalysis Using Models of Heterogeneous Cover Signals*. PhD thesis, Faculty of Comp. Sci., TU Dresden, Germany, 2008.

[2] T. Filler and J. Fridrich. Gibbs construction in steganography. *IEEE TIFS*, 5(4):705–720, 2010.

[3] T. Filler and J. Fridrich. Design of adaptive stegano-graphic schemes for digital images. In *Proc. SPIE, Elec. Img., Media Watermarking, Sec. and Forensics of Multimedia XIII*, volume 7880, pages OF 1–14, 2011.

[4] T. Filler, J. Judas, and J. Fridrich. Minimizing additive distortion in steganography using syndrome-trellis codes. *IEEE TIFS*, 6(3):920–935, September 2011.

[5] T. Filler, T. Pevný, and P. Bas. BOSS (Break Our Steganography System). http://www.agents.cz/boss, July 2010.

[6] J. Fridrich and J. Kodovský. Rich models for steganalysis of digital images. *IEEE TIFS*, 7(3):868–882, 2011.

[7] J. Fridrich, J. Kodovský, M. Goljan, and V. Holub. Steganalysis of content-adaptive steganography in spatial domain. In *Information Hiding, 13th Int. Conf.*, volume 6958 of *Springer LNCS*, pages 102–117, 2011.

[8] G. Gül and F. Kurugollu. A new methodology in steganalysis : Breaking highly undetactable steganograpy (HUGO). In *Information Hiding, 13th Int. Conf.*, volume 6958 of *Springer LNCS*, pages 71–84, 2011.

[9] Y. Kim, Z. Duric, and D. Richards. Modified matrix encoding technique for minimal distortion steganography. In *Information Hiding, 8th Int. Workshop*, volume 4437 of *Springer LNCS*, pages 314–327, 2006.

[10] J. Kodovský and J. Fridrich. On completeness of feature spaces in blind steganalysis. In *Proc. of the 10th ACM MM&Sec Workshop*, pages 123–132, 2008.

[11] J. Kodovský, J. Fridrich, and V. Holub. On dangers of overtraining steganography to incomplete cover model. In *Proc. of the 13th ACM MM&Sec Workshop*, pages 69–76, 2011.

[12] J. Kodovský, J. Fridrich, and V. Holub. Ensemble classifiers for steganalysis of digital media. *IEEE TIFS*, 7(2):432–444, 2012.

[13] W. Luo, F. Huang, and J. Huang. Edge adaptive image steganography based on LSB matching revisited. *IEEE TIFS*, 5(2):201–214, June 2010.

[14] T. Pevný, T. Filler, and P. Bas. Using high-dimensional image models to perform highly undetectable steganography. In *Information Hiding, 12th Int. Conf.*, volume 6387 of *Springer LNCS*, pages 161–177, 2010.

[15] V. Sachnev, H. J. Kim, and R. Zhang. Less detectable JPEG steganography method based on heuristic optimization and BCH syndrome coding. In *Proc. of the 11th ACM MM&Sec Workshop*, pages 131–140, 2009.

[16] C. Wang and J. Ni. An efficient JPEG steganographic scheme based on the block–entropy of DCT coefficents. In *Proc. of IEEE ICASSP*, Kyoto, Japan, 2012.