

# A Theoretical Framework for Data-Hiding in Digital and Printed Text Documents

R. Villán, S. Voloshynovskiy, F. Deguillaume, Y. Rytsar,  
O. Koval, E. Topak, E. Rivera, and T. Pun

Computer Vision and Multimedia Laboratory - University of Geneva,  
24, rue du Général-Dufour - 1211 Geneva 4, Switzerland  
svolos@cui.unige.ch

In this work, we consider the text data-hiding problem as a particular instance of the well-known Gel'fand-Pinsker problem [1]. The text, where some message  $m \in \mathcal{M}$  is to be hidden, is represented by  $\mathbf{x}$  and called cover text. Each component  $x_i$ ,  $i = 1, 2, \dots, N$ , of  $\mathbf{x}$  represents one character from this text. Here, we define a character as an element from a given language alphabet (e.g. the latin alphabet  $\{A, B, \dots, Z\}$ ). To be more precise, we conceive each character  $x_i$  as a data structure consisting of multiple component fields (features): *name*, *shape*, *position*, *orientation*, *size*, *color*, etc.

Assuming the knowledge of the conditional probability distribution  $p(u|x)$ ,  $|\mathcal{M}||\mathcal{J}|$  codewords  $\mathbf{u}$  are generated independently at random and located into  $|\mathcal{M}|$  bins, each of them with  $|\mathcal{J}|$  codewords. Once generated, the codebook is revealed to both the encoder and the decoder. Given  $m$  to be communicated, the encoder produces the watermark  $\mathbf{w}$  by finding first a jointly strongly typical pair  $(\mathbf{x}, \mathbf{u}(m, j))$ , where  $\mathbf{u}(m, j)$  is the  $j$ -th codeword inside the bin corresponding to  $m$ , and then, by using a deterministic mapping  $\mathbf{w} = \varphi^N(\mathbf{x}, \mathbf{u})$ . The influence of the channel  $p(v|w, x)$  is divided in two stages. In the first stage,  $\mathbf{w}$  and  $\mathbf{x}$  are combined via a deterministic mapping  $\psi^N(\mathbf{w}, \mathbf{x})$  to give the stego text  $\mathbf{y}$ . In the second stage,  $\mathbf{y}$  may suffer from some intentional or unintentional distortions. We denote by  $\mathbf{v}$  the resulting distorted version of  $\mathbf{y}$ . Finally,  $\mathbf{v}$  is fed to the decoder, which tries to obtain an estimate  $\hat{m}$  of message  $m$  by using the jointly strongly typical decoding rule.

As a particular example of the Gel'fand-Pinsker scheme, let us consider the Scalar Costa Scheme (SCS) [2] where the stego text  $Y$  is obtained as  $Y = W + X = \alpha' Q_m(X) + (1 - \alpha')X$ , where  $Q_m(\cdot)$  is a scalar quantizer corresponding to  $m$  and  $\alpha'$  is a compensation parameter. For a practical implementation based on the SCS, we only need to select a character feature (e.g. color), and use it as the cover character  $X$ . We show in Fig. 1 the resulting SCS codebook and an illustration of how to use it for text data-hiding.

Based on the above framework, we propose two new methods for text data-hiding: *color quantization* and *halftone quantization*. The exploited character features are, respectively, *color* and *halftone pattern* (see Fig. 2). The main idea of these methods is to quantize the character feature in such a manner that the human visual system is not able to distinguish between the original and quantized characters, but it is still possible to do it by a specialized reader, e.g. a high dynamic range and/or high resolution scanner in the case of printed

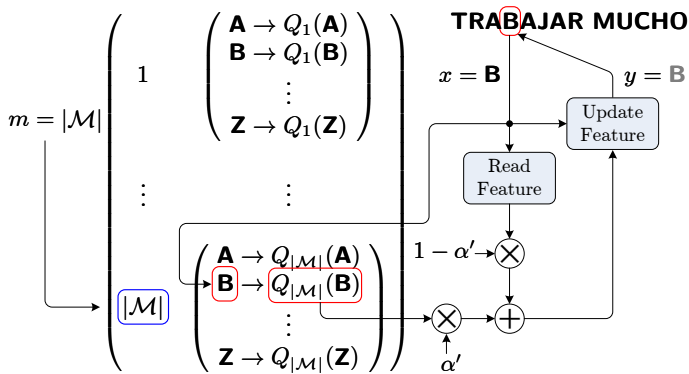


Fig. 1. SCS text data-hiding

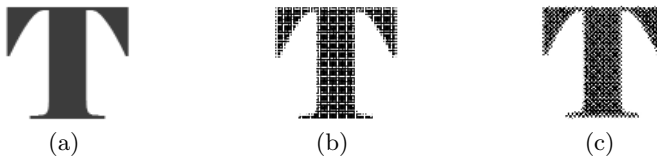


Fig. 2. Halftone quantization: (a) original character; (b) marked character for  $m = 0$ , screen angle =  $0^\circ$ ; (c) marked character for  $m = 1$ , screen angle =  $45^\circ$

documents. In particular, we show that the color quantization method works both for digital and printed documents, has high information embedding rate, is perceptually invisible, and is fully automatable.

## Acknowledgment

This work was partially supported by the SNSF professorship grant no. PP002-68653/1, the IM2 project, and the European Commission through the programs IST-2002-507932 ECRYPT and FP6-507609 SIMILAR. The information in this document reflects only the author's views, is provided as is and no guarantee or warranty is given that the information is fit for any particular purpose. The user thereof uses the information at its sole risk and liability.

## References

1. Gel'fand, S., Pinsker, M.: Coding for channel with random parameters. *Problems of Control and Information Theory* **9** (1980) 19–31
2. Eggers, J., Su, J., Girod, B.: A blind watermarking scheme based on structured codebooks. In: *Secure Images and Image Authentication*, IEE Colloquium, London, UK (2000) 4/1–4/6