

---

# Embedding Secret Data in Html Web Page

Sandipan Dey<sup>1</sup>, Hameed Al-Qaheri<sup>2</sup>, and Sugata Sanyal<sup>3</sup>

<sup>1</sup> Cogniant Technology Solutions [sandipan.dey@gmail.com](mailto:sandipan.dey@gmail.com)

<sup>2</sup> Department of Quantitative Methods and Information Systems College of Business Administration Kuwait University [alqaheri@cba.edu.kw](mailto:alqaheri@cba.edu.kw)

<sup>3</sup> School of Technology and Computer Science Tata Institute of Fundamental Research, Homi Bhabha Road, Mumbai - 400005, India [sanyal@tifr.res.in](mailto:sanyal@tifr.res.in)

**Summary.** In this paper, we suggest a novel data hiding technique in an Html Web page. Html Tags are case insensitive and hence an alphabet in lowercase and one in uppercase present inside an html tag are interpreted in the same manner by the browser, i.e., change in case in an web page is imperceptible to the browser. We basically exploit this redundancy and use it to embed secret data inside an web page, with no changes visible to the user of the web page, so that he can not even suspect about the data hiding. The embedded data can be recovered by viewing the source of the html page. This technique can easily be extended to embed secret message inside any piece of source-code where the standard interpreter of that language is case-insensitive.

## 1 Introduction

Some techniques for hiding data in executables are already proposed (e.g., Shin et al [4]). In this paper we introduce a very simple technique to hide secret message bits inside source codes as well. We describe our steganographic technique by hiding inside html source as cover text, but this can be easily extended to any case-insensitive language source codes. Html Tags are basically directives to the browser and they carry information regarding how to structure and display the data on a web page. They are not case sensitive, so tags in either case (or mixed case) are interpreted by the browser in the same manner (e.g., "< head >" and "< HEAD >" refers to the same thing). Hence, there is a redundancy and we can exploit this redundancy. To embed secret message bits into html, if the cases of the tag alphabets in html cover text are accordingly manipulated, then this tampering of the cover text will be ignored by the browser and hence it will be imperceptible to the user, since there will not be any visible difference in the web page, hence there will not be any suspect for it as well. Also, when the web page is displayed in the browser, only the text contents are displayed, not the tags (those can only be seen when the user does 'view source'). Hence, the secret messages will be kind of hidden to user.

Both redundancy and imperceptibility conditions for data hiding are met, we use these to embed data in html text. If we do not tamper the html text data that is to be displayed by the browser as web page (this html cover text is analogical

to the cover image, when thought in terms of steganographic techniques in images [1, 2, 3]), the user will not even suspect about hidden data in text. We shall only change the case of every character within these Html tags (elements) in accordance with the secret message bits that we want to embed inside the html web page. If we think of the browser interpreter as a function,  $f_B : \Sigma^* \rightarrow \Sigma^*$  we see that it is non-injective, i.e., not one to one, since  $f_B(x) = f_B(y)$  whenever  $x \in \{ 'A' \dots 'Z' \}$ ,  $y \in \{ 'a' \dots 'z' \}$  and  $Uppercase(y) = x$ . The extraction process of the embedded message will also be very simple, one needs to just do 'view source' and observe the case-patterns of the text within tags and can readily extract the secret message (and see the unseen), while the others will not know anything.

The length (in bits) of the secret message to be embedded will be upperlimited by the sum of size of text inside html tags (here we don't consider attribute values for data embedding. In case we consider attribute values for data embedding, we need to be more careful, since for some tags we should think of case-sensitivity, e.g. `<A HREF="link.html">`, since link file name may be casesensitive on some systems, whereas, attributes such as `<h2 align="center">` is safe). If less numbers of bits to be embedded, we can embed the information inside Header Tag specifying the length of embedded data (e.g. `'<Header 25 >'` if the length of secret data to be embedded is 25 bits) that will not be shown in the browser (optionally we can encrypt this integer value with some private key). In order to guarantee robustness of this very simple algorithm one may use some simple encryption on the data to be embedded.

## 2 The Algorithm for Embedding

The algorithm for embedding the secret message inside the html cover text is very simple and straight-forward. First, we need to separate out the characters from the cover text that will be candidates for embedding, these are the caseinsensitive text characters inside Html tags. Figure 2 shows a very simplified automata for this purpose.

We define the following functions before describing the algorithm:

- $l : \Sigma^* \rightarrow \Sigma^*$  is defined by,

$$l(c) = \begin{cases} ToLower(c) & c \in \{ 'A' \dots 'Z' \} \\ c & \text{otherwise} \end{cases} \text{ were } ToLower(c) = c + 32$$

- Similarly,  $u : \Sigma^* \rightarrow \Sigma^*$  is defined by,

$$u(c) = \begin{cases} ToUpper(c) & c \in \{ 'a' \dots 'z' \} \\ c & \text{otherwise} \end{cases} \text{ were } ToUpper(c) = c - 32$$

Here the ascii value of 'A' is 65 and that of 'a' is 97, with a difference of 32.

It's easy to see that if the domain  $\Sigma^* = \{ 'a' \dots 'z' \} \cup \{ 'A' \dots 'Z' \}$ , then  $l : \{ 'A' \dots 'Z' \} \rightarrow \{ 'a' \dots 'z' \}$  and  $u : \{ 'a' \dots 'z' \} \rightarrow \{ 'A' \dots 'Z' \}$ , implies that  $l(.) = u(.) = \Sigma^* - l(.)$ .

Now, we want to embed secret data bits  $b_1, b_2 \dots b_k$  inside the case-insensitive text inside the Html Tags. If  $c_1 c_2 \dots c_n$  denotes the sequence of characters inside the html tags in cover text (input html). A character  $c_i$  is a candidate for hiding a

secret message bit iff it is an alphabet. If we want to hide the  $j^{\text{th}}$  secret message bit  $b_j$  inside the cover text character  $c_i$ , the corresponding stego-text will be defined by the following function  $f_{stego}$ :  $\forall c_i \in \{ 'a' \dots 'z' \} \cup \{ 'A' \dots 'Z' \}$ , i.e. if  $IsAlphabet(c_i)$  is true,  $f_{stego}(c_i) = \begin{cases} l(c_i) & b_j = 0 \\ u(c_i) & b_j = 1 \end{cases}$  Hence, we have the following:

$$c_i \in \{ 'a' \dots 'z' \} \cup \{ 'A' \dots 'Z' \} \Rightarrow f_{stego}(c_i) = l(c_i) \cdot \overline{b_j} + u(c_i) \cdot b_j, \forall i \quad (1)$$

Number of bits ( $k$ ) of the secret message embedded into the html cover text must also be embedded inside the html (e.g., in Header element). The figure 1 and the algorithm 1 together explain this embedding algorithm.

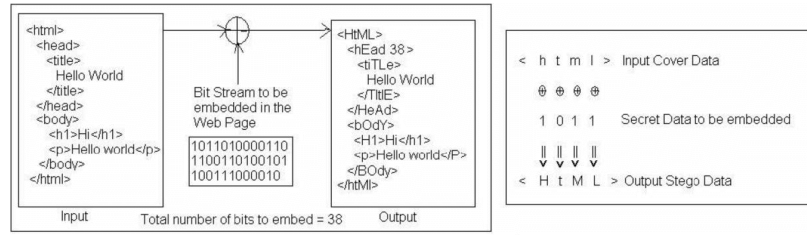


Fig. 1: Illustration of how the data hiding works

### 3 The Algorithm for Extraction

---

#### Algorithm 1 Embedding Algorithm

---

Search for all the html tags present in the html cover text and extract all the characters  $c_1 c_2 \dots c_n$  from inside those tags using the DFA described in the Fig.2. Embed the secret message length  $k$  inside html header in the stego text.

$j \leftarrow 0$ .

**for**  $c_i \in HTMLTAGS$ ,  $i = 1 \dots n$  **do**

**if**  $c_i \in \{ 'a' \dots 'z' \} \cup \{ 'A' \dots 'Z' \}$  **then**

$f_{stego}(c_i) = l(c_i) \cdot \overline{b_j} + u(c_i) \cdot b_j$ .

$j \leftarrow j + 1$ .

**else**

$f_{stego}(c_i) = c_i$ .

**end if**

**if**  $j == k$  **then**

        break.

**end if**

**end for**

---

The algorithm for extraction of the secret message bits will be even more simple. Like embedding process, we must first separate out the candidate text (text within tags) that were chosen for embedding secret message bits. Also, we must extract the number of bits ( $k$ ) embedded into this page (e.g., from the header element). One has to use 'view source' to find out the stego-text.

Now, we have  $d_i = f_{stego}(c_i)$ ,  $\forall i \in \{1, 2, \dots, n\}$ . If  $d_i \in \{'a' \dots 'z'\} \cup \{'A' \dots 'Z'\}$  i.e., an alphabet, then only it is a candidate for decoding and to extract  $b_i$  from  $d_i$ , we use the following logic:  $b_i = \begin{cases} 0 & d_i \in \{'a' \dots 'z'\} \\ 1 & d_i \in \{'A' \dots 'Z'\} \end{cases}$ . Repeat the above algorithm  $\forall i < k$ , to extract all the hidden bits. The figure 2 and the algorithm 2 together explain this embedding algorithm.

---

**Algorithm 2** Extraction Algorithm

---

Search for all the html tags present in the html stego text and extract all the characters  $d_1 d_2 \dots d_n$  from inside those tags using the DFA described in the Fig.2. Extract the secret message length  $k$  from inside html header in the stego text.

```

j ← 0.
for  $d_i \in HTMLTAGS$ ,  $i = 1 \dots n$  do
  if  $d_i \in \{'a' \dots 'z'\}$  then
     $b_j = 0$ .
     $j \leftarrow j + 1$ .
  else
    if  $d_i \in \{'A' \dots 'Z'\}$  then
       $b_j = 1$ .
       $j \leftarrow j + 1$ .
    end if
  end if
  if  $j == k$  then
    break.
  end if
end for

```

---

Figures 3, 4 and 5 show an example of how our method works, while Figure 6 shows the comparison of the histogram of the cover html and stego html in terms of the (ascii) character frequencies. Classical image hiding techniques like LSB data hiding technique always introduce some (visible) distortion [5, 10] in the stego image (that can be reduced using techniques [6, 7, 8, 9]), but our data hiding technique in html is novel in the sense that it introduces no visible distortion in stego text at all.

## 4 Conclusions

In this paper we presented an algorithm for hiding data in html text. This technique can be extended to any case-insensitive language and data can

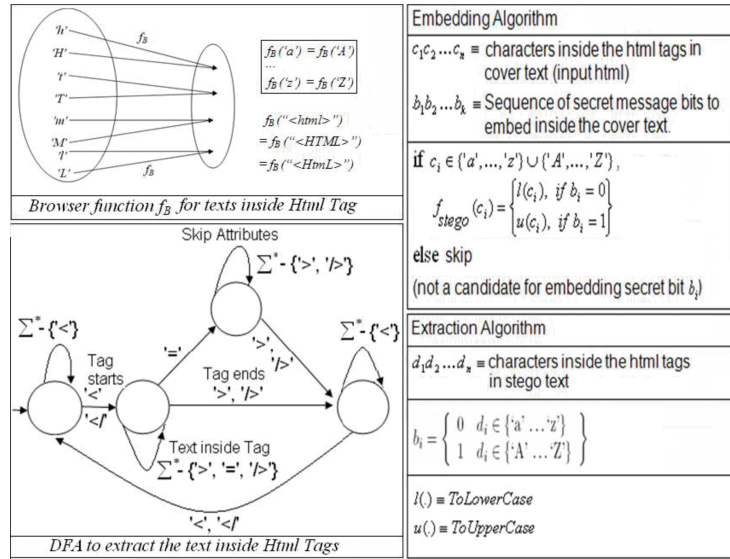


Fig. 2: Basic block-diagram for html data-hiding technique

be embedded in the similar manner, e.g., we can embed secret message bits even in source codes written in languages like basic or pascal or in the case-insensitive sections (e.g. comments) in C like case-sensitive languages. Data hiding methods in images results distorted stego-images, but the html data hiding technique does not create any sort of visible distortion in the stego html text.

## References

1. Neil F. Johnson and Sushil Jajodia. Steganography: Seeing the Unseen IEEE Computer, February 1998, pp. 26-34.
2. W. Bender et al., Techniques for Data Hiding, Systems J., Vol. 35, Nos. 3 and 4, 1996, pp. 313-336.
3. B. Pfitzmann, Information Hiding Terminology, Proc. First International Workshop Information Hiding, Lecture Notes in Computer Science No. 1,174, Springer-Verlag, Berlin, 1996, pp. 347-356.
4. DaeMin Shin, Yeog Kim, KeunDuck Byun, SangJin Lee, Data Hiding in Windows Executable Files, Center for Information Security Technologies (CIST), Korea University, Seoul, Republic of Korea
5. C. Kurak, J. McHugh, A Cautionary Note On Image Downgrading, Proc. IEEE Eighth Annual Computer Security Applications Conf., IEEE Press, Piscataway, N.J., 1992, pp. 153-159.
6. F. Battisti, M. Carli, A. Neri, K. Egiazarian, A Generalized Fibonacci LSB Data Hiding Technique, 3rd International Conference on Computers and De-

[Cover Html Source](#)

Fig. 3: Cover Html source before embedding the secret message

- vices for Communication (CODEC- 06) TEA, Institute of Radio Physics and Electronics, University of Calcutta, December 18-20, 2006.
7. Sandipan Dey, Ajith Abraham, Bijoy Bandyopadhyay and Sugata Sanyal, Data Hiding Techniques Using Prime and Natural Numbers, Journal of Digital Information Management, Volume 6, 2008.
8. Sandipan Dey, Ajith Abraham and Sugata Sanyal An LSB Data Hiding Technique Using Natural Numbers, IEEE Third International Conference on Intelligent Information Hiding and Multimedia Signal Processing, IIHMSP 2007, Nov 26-28, 2007, Kaohsiung City, Taiwan, IEEE Computer Society press, USA, ISBN 0-7695-2994-1, pp. 473-476, 2007.
9. Sandipan Dey, Ajith Abraham and Sugata Sanyal An LSB Data Hiding Technique Using Prime Numbers, Third International Symposium on Information Assurance and Security, August 29-31, 2007, Manchester, United Kingdom, IEEE Computer Society press, USA, ISBN 0-7695-2876-7, pp. 101-106, 2007.
10. R. Z. Wang, C. F. Lin and I. C. Lin, Image Hiding by LSB substitution and genetic algorithm, Pattern Recognition, Vol. 34, No. 3, pp. 671-683, 2001.



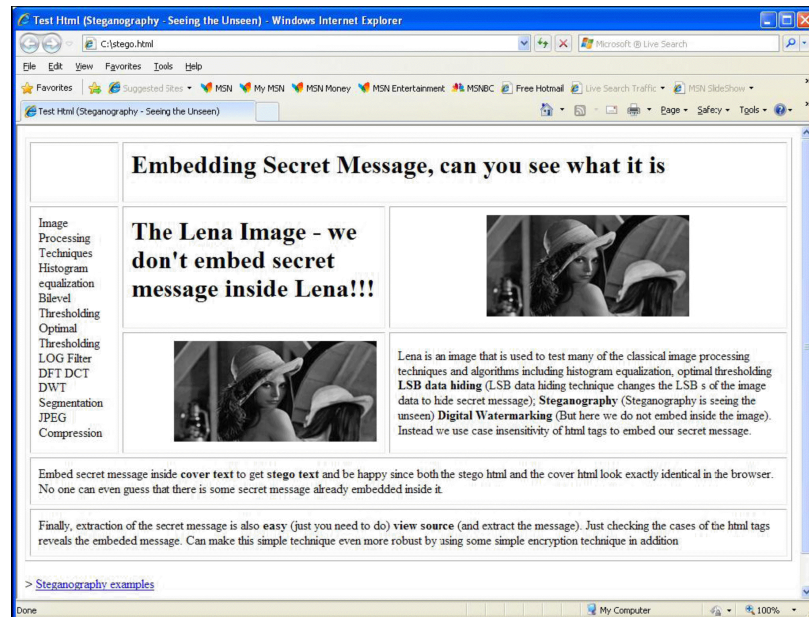


Fig. 5: Cover & Stego Html

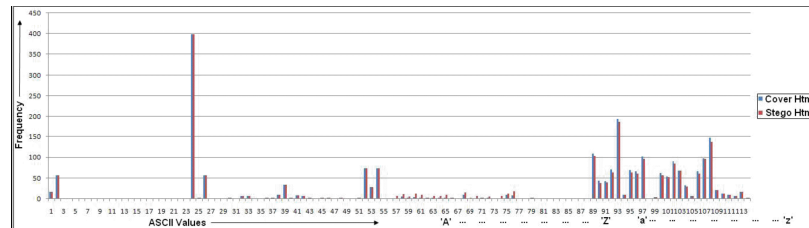


Fig. 6: Cover vs Stego Html Histogram