

New Steganographic Techniques for the OOXML File Format

Aniello Castiglione^{1,*}, Bonaventura D'Alessio¹,
Alfredo De Santis¹, and Francesco Palmieri²

¹ Dipartimento di Informatica “*R. M. Capocelli*”
Università degli Studi di Salerno, I-84084 Fisciano (SA), Italy
Tel.:+39089969594; Fax: +39089969594

castiglione@{ieee,acm}.org, bdalessio@dia.unisa.it, ads@dia.unisa.it

² Dipartimento di Ingegneria dell'Informazione
Seconda Università degli Studi di Napoli, I-81031 Aversa (NA), Italy
francesco.palmieri@unina.it

Abstract. The simplest container of digital information is the file and among the vast array of files currently available, MS-Office files are the most widely used. The “Microsoft Compound Document File Format” (MCDF) has often been used to host secret information. The new format created by Microsoft, first used with MS-Office 2007, makes use of a new standard, the “Office Open XML Formats” (OOXML). The benefits include that the new format introduces the OOXML format, which lowers the risk of information leakage, as well as the use of MS-Office files as containers for steganography.

This work presents some new methods of embedding information into the OOXML file format which can be extremely useful when using MS-Office documents in steganography. The authors highlight how the new methods introduced in this paper can also be used in many other scenarios, not only in MS-Office documents. An evaluation of the limits of the proposed methods is carried out by comparing them against the tool introduced by Microsoft to sanitize MS-Office files. The methods presented can be combined in order to extend the amount of data to be hidden in a single cover file.

Keywords: Steganography, OOXML Format, Stegosystem, Document Steganography, Microsoft Office Document, Information Hiding, Document Metadata, Covert Channel.

1 Introduction

The MS-Office suite is, without a doubt, the most widely used word-processing tool for preparing and writing documents, spreadsheets and presentations [13]. Therefore, the possibility to hide information inside them is a challenge that probably concerns many different parties. Starting with the 2007 version (MS-Office 2007), Microsoft has completely changed the format of its files increasing,

* Corresponding author

among other things, the level of security and thus making it more difficult to hide information inside them. In fact, it has gone from using the old binary format to the new OOXML [5], which uses XML files. In addition to guarantee a significantly high level of “privacy and security”, it has also introduced the feature *Document Inspector*, which makes it possible to quickly identify and remove any sensitive, hidden and personal information. It is therefore evident that the old methodologies of Information Hiding that exploit the characteristics of the binary files of MS-Office are no longer applicable to the new XML structures. However, the steganography techniques that take advantage of the functions offered by the Microsoft suite([7], [8], [9], [10]), are still valid, and therefore independent from the version used. The new format offers new perspectives, as proposed by Garfinkel et al. [6] as well as Park et al. [15]. Both authors describe methodologies which use characteristics that do not conform to the OOXML standard and therefore can be characterized by searching for abnormal content type that is not described in the OOXML specifications inside the file.

This study proposes and analyzes four new steganography techniques for MS-Office files, with only the first not taking advantage of characteristics that do not conform to the OOXML standard.

The remaining of this paper is structured as follows. Section 2 introduces the OOXML standard and the features of the Document Inspector. Section 3 discusses the methodology that takes advantage of the possibility to use different compression algorithms in generating MS-Office files. Section 4 highlights how it is possible to hide data in the values of the attribute that specifies a unique identifier used to track the editing session (revision identifier). In Section 5 a methodology, that uses images not visualized by MS-Office, but present in the file, is analyzed in order to hide information. Section 6 illustrates how the macro of MS-Office can be used to hide information. In Section 7 the methodologies are compared, verifying the overhead introduced as well as the resulting behavior of save actions.

2 The OOXML Format

Starting with the 2007 version, Microsoft has adopted the OOXML format based on XML (XML-based file format). In fact, Microsoft has begun the transition from the old logic, that saw the generation of a binary file format, to a new one that uses XML files. The Extensible Markup Language (XML) is used for the representation of structured data and documents. It is a markup language and, thus, composed of instructions, defined as tags or markers. Therefore, in XML a document is described, in form and content, by a sequence of elements. Every element is defined by a *tag* or a pair *start-tag/end-tag*, which can have one or more attributes. These attributes define the properties of the elements in terms of values. The OOXML format is based on the principle that even a third party, without necessarily owning product rights, can extract and relocate the contents of the MS-Office file by only using standard transformation methods. This is possible because XML text is clearly written and therefore visible and

modifiable with any text editor. Moreover, OLE attachments are present in the source file format and therefore can be visualized with any compatible viewer.

Distinguishing documents produced in this new format is easy due to the file extensions being characterized by an “*x*” at the end, with the file Word, Excel and PowerPoint respectively being *.docx*, *.xlsx*, *.pptx*. An additional feature is that a macro is not activated unless specified by the user. In this case, the extension of the files changes by adding “*m*” rather than “*x*” and thus become *.docm*, *.xlsm*, *.pptm*. The new structure of an OOXML file, which is based on the ECMA-376 standard [3], uses a container, a ZIP file, inside of which there are a series of files, mostly XML, and are opportunely organized into folders, that describe both the content as well as the properties and relationships of them. It is highly likely that the ZIP standard was chosen because it is the most commercially well-known, in addition to having characteristics of flexibility and modularity that allow for any eventual expansions in future functionalities [16]. There are three types of files stored in the “container”, that can be common to all the applications of MS-Office or specific for each one (Word, Excel, PowerPoint):

- XML files, that describe application data, metadata, and even customer data, stored inside the container file;
- non-XML files, may also be included within the container, including such parts as binary files representing images or OLE objects embedded in the document;
- relationship parts that specify the relationships between the parts; this design provides the structure for an MS-Office file.

For example, analyzing a simple Word document, the structure [4] of the folders and files in a ZIP container will be like that shown in Fig. 1.

Therefore, beginning from version 2007, the MS-Office documents:

- are files based on the ZIP standard;
- contain XML files;
- have common characteristics and formats to those of generic MS-Office files (character format, cell properties, collaborative document, etc.);
- may contain OLE objects (images, audio files, etc.);
- conform to the ECMA-376 standard, opportunely customized.

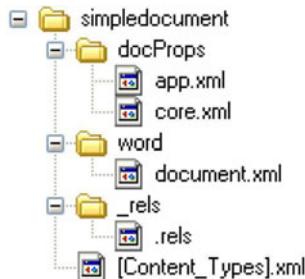


Fig. 1. Structure of a simple Word document

Another key concept related to the OOXML format is the modularity, either inside the files or among files, which allows for either the easy addition of new elements or the removal of old ones. For example, the addition of a new JPEG image inside a Word file could be simply performed by:

- copying the file with the *.jpg* extension in the folder named *media* within the ZIP container;
- adding a group of elements in the *document.xml* file (it contains the XML markup that defines the contents of the document) in order to describe the insertion methods within the page;
- adding, in several files of the relationship, some XML lines which declare the use of an image.

The OOXML format gives new opportunities to the community, as indicated by Microsoft [5]. In fact with the new standard:

- it is possible to show only the text of the document. If the file is a Word document, for example, only the file *document.xml* will be analyzed without necessarily opening all the files which contain the remaining information about the document;
- the files are compressed, and consequently are shorter and easy to manage;
- it is simpler to scan for viruses or malicious contents thanks to its textual form instead of the old binary format;
- the new format does not allow to have macro inside it, thus guaranteeing a satisfactory level of security;
- if some of the files in the ZIP container are damaged, the integrity of the entire document could be preserved, and in some cases the main document could be reconstructed starting from the remaining “untouched” files.

MS-Office 2010, also known as Office 14, maintains formats and interfaces that are similar to the 2007 version. The substantial difference between the two suites is that MS-Office 2010 is much more web-oriented than the previous one. The new suite, for example, sends the user an alert message when transmitting sensitive information via e-mail. It is also able to translate documents and deal with different languages, as well as transform presentations into clips. It makes possible to present a PowerPoint “slideshow” to users connected to the Internet. In [11] Microsoft analyzes, describing some of their characteristics, all the new features introduced in the new version, highlighting the updated parts with respect to the old version.

The management flexibility offered by the new OOXML format has obvious implications when dealing with security. On one hand, the clear-text offers the seeming impossibility to hide information. While, on the other, it offers the possibility to malicious parties to read its content and eventually freely manipulate it. It is also well-known that MS-Office files contain data that can reveal unwanted personal information, such as people who have collaborated in the writing of the document, network parameters, as well as devices on which it has been edited. In current literature, there are several papers which describe how to extract and

reconstruct different types of information from such documents. Castiglione et al. [1] introduced a steganography system which can be applied to all versions before MS-Office 2007. Furthermore, authors analyzed the information leakage issue raised prior of MS-Office 2007 documents.

In order to guarantee a higher level of security and privacy, Microsoft (starting from MS-Office 2007 for Windows) have introduced the feature called Document Inspector which makes it possible to find and remove, quickly, personal, sensitive and hidden information. More details on the Document Inspector can be found in [12].

3 Data Hiding by Different Compression Algorithm of ZIP

Taking advantage of the characteristic that OOXML standard produces compressed files, it is possible to hide information inside a ZIP structure without taking into account that the same file will be interpreted by MS-Office as a document produced by its own application. The ZIP format is a data compression and archive format. Data compression is carried out using the DeflateS format [2], which is set as default, with it being possible to set a different compression algorithm. For example, by using WinZip (ver. 14.5 with the command-line add-on ver. 3.2) it is possible to choose one of the compression algorithm indicated in Table 1.

Therefore, by inserting in the command

```
wzip [options] zipfile [@listafilename] [files...]
```

one of the options indicated in Table 1, the desired algorithm compression will be applied. It is worth noting that, in a ZIP container, all the files contained can be compressed with a different algorithm. In MS-Office files, that are ZIP containers, it is possible to set various compression algorithms.

Not all the algorithms listed in Table 1 are correctly interpreted by MS-Office. In fact, after some tests, it has been possible to ascertain that only the 5 algorithms present in Table 2 are supported by MS-Office. Initially, the tests has

Table 1. Compression options in the ZIP format

Algorithm	Acronym	Option
maximum (<i>PPMd</i>)	<i>PPDM</i>	<i>ep</i>
maximum (<i>LZMA</i>)	<i>LZMA</i>	<i>el</i>
maximum (<i>bzip2</i>)	<i>BZIPPED</i>	<i>eb</i>
maximum (<i>enhanced deflate</i>)	<i>EnhDefl</i>	<i>ee</i>
maximum (<i>portable</i>)	<i>DeflateX</i>	<i>ex</i>
normal	<i>DeflateN</i>	<i>en</i>
fast	<i>DeflateF</i>	<i>ef</i>
super fast	<i>DeflateS</i>	<i>es</i>
best method for each file (based on the file type)		<i>ez</i>
no compression	<i>Stored</i>	<i>e0</i>

Table 2. Association character-algorithms

Algorithm	Option	Char
DeflatF	<i>ef</i>	0
DeflatN	<i>en</i>	1
DeflatX	<i>ex</i>	2
DeflateS	<i>es</i>	3
Stored	<i>e0</i>	4

been performed on a .docx file, which has been compressed by using the different compression algorithms. It has been determined that both MS-Office 2007 and MS-Office 2010 do not correctly handle file compressed with the following compression switches: *eb*, *ee*, *el*, *ep*, *ez*. In such a case, it is shown an error message stating that the ZIP format is not supported. MS-Office uses by default the compression algorithm named *DeflateS*.

The proposed steganographic technique considers different compression algorithms as different parameters of source encoding. More precisely:

- hidden data is codified with an alphabet of 5 elements, the 5 different values that indicate the compression algorithm used;
- the codes obtained through the previous point are hidden in ZIP files associating a character to every file present in the container;
- the compression algorithm applied to the single file corresponds to the value of the character to be hidden.

Example 1. Consider the binary string $(1010101101111111001000100001)_2$ to be hidden in a Word document which has just been created and has no characters. This document is made up of 12 files, as listed in the first column of Table 3.

Table 3. Decoding table

File	Algorithm	Char
<i>[ContentTypes].xml</i>	DeflatS	3
<i>\docProps\app.xml</i>	DeflatS	3
<i>\docProps\core.xml</i>	DeflatX	2
<i>\word\document.xml</i>	DeflatF	0
<i>\word\fontTable.xml</i>	DeflatN	1
<i>\word\settings.xml</i>	DeflatS	3
<i>\word\styles.xml</i>	Stored	4
<i>\word\stylesWithEffects.xml</i>	DeflatS	3
<i>\word\webSettings.xml</i>	DeflatX	2
<i>\word\theme\theme1.xml</i>	Stored	4
<i>\word_rels\document.xml.rels</i>	DeflatN	1
<i>_rels_rels</i>	DeflatS	3

The files are listed in alphabetical order in relation to their “absolute” name (comprehensive of the path). Thus, there is an univocal sequence on which it codifies or decodes. In order to hide the binary string, it has to be first converted into a number in base 5. The base 5 representation of the number $(101010110111111001000100001)_2$ is a string of 12 numbers: $(332013432413)_5$. It is assumed that the values indicated in Table 2 can be associated to the various compression algorithms. In order to obtain the stego-text, every file will be simply compressed with the corresponding algorithm associated to the character to be hidden (see Table 3). \square

If the MS-Office file contains M files, the proposed technique allows to hide

$$\log_2 5^M = M \cdot \log_2 5 \cong M \cdot 2.32$$

bits of information. M is at least 12, but usually is greater.

4 Data Hiding by the Revision Identifier Value

The second proposed method of hiding information in MS-Office documents, which is only applicable to Word files, is to use the value of several attributes that are in XML. The revision identifier *rsid* is a sequence of 8 characters which specifies a unique identifier used to track the editing session. An editing session is defined as the period of editing which takes place between any two subsequent save actions. The *rsid*, as an attribute of an XML element, gives information on the part of code contained in the same element. The types of revision identifier, usable in the OOXML standard, are listed in the specifications of the ECMA-376. These attributes, defined as the *ST_LongHexNumber* simple type, are strings of 8 hexadecimal characters:

$$(x_0x_1x_2x_3x_4x_5x_6x_7) : x_i \in \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, A, B, C, D, E, F\}$$

All the revision identifier attributes, present with the same value in a document, indicate that the code in the element has been modified during the same editing session.

An example element which contains 3 *rsid* attributes is:

```
<w:p w:rsidR="000E634E" w:rsidRDefault="008C3D74" w:rsidP="00463DF8">
```

It is worth noting that there are three sequences of 8 characters, that represent the unique identifier associated to the attributes: *rsidR*, *rsidRDefault* and *rsidP* (see pp. 243-244 of the ECMA-376 specifications [3]).

The methodology proposed in this section consists of replacing the values of the *rsid* attributes with the data to be hidden, codified in hexadecimal. Thus, if T is the number of occurrences of these attributes in the MS-Office files, the maximum number of bits that can be hidden will be:

$$\log_2 16^{T \cdot 8} = 32 \cdot T$$

due to every attribute being composed of 8 hexadecimal characters. If the information to be hidden exceeds the maximum number of bits that can be contained in the MS-Office document, it is possible to add to the XML file further elements with *rsid* attributes. Furthermore, one more trick is required to avoid the detection of hidden data by a stego-analysis inspection. MS-Office records in the file `setting.xml` all the *rsid* values that have been used in different versions of the file `document.xml`. To perform such an activity, MS-Office uses the XML element `<w:rsid w:val="002A31DF">`. Consequently, when, to hide information, it is used the methodology presented in this section, after having modified the *rsid* values in the file `document.xml`, it is necessary to insert the same values even in the file `setting.xml`. In fact, the presence of *rsid* values in the file `document.xml` which are not present in the file `setting.xml` is a strange situation that could raise suspicion.

Among the various functionalities available in MS-Office, there is the possibility to track the changes of a document. By using such feature, MS-Office keeps track of all the modifications performed in a document (deleted, inserted or modified text), of the date when they have been made and of the user who has carried out such modifications. Those information, even though can be partially reconstructed by the analysis of the *rsids*, are traced by using two XML elements. Such elements, delimited by a pair of *start-tag* and *end-tag*, are different if used to track a deletion (with the tag `<w:del . . . > </w:del>`) or an insertion (with the tag `<w:ins . . . > </w:ins>`).

This element has the following 3 attributes: identification code (*id*), author who modified the document (*author*) as well as time and date in which the change (*date*) occurred (this is an optional attribute). Consequently, all the modifications performed by the same author within the same editing session will be placed in the XML file between the *start-tag* and *end-tag* of the “change-tracking” element.

For example, if the user PCCLIENT would have deleted the text “one” at 09:23:00 GMT of October 11, 2010, the code excerpt will be like:

```
<w:del w:id="0" w:author="PCCLIENT" w:date="2010-10-11T09:23:00Z">
  <w:r w:rsidRPr="00111111" w:rsidDel="00333333">
    <w:rPr>
      <w:lang w:val="en-US"/>
    </w:rPr>
    <w:delText xml:space="preserve">one</w:delText>
  </w:r>
</w:del>
```

That being stated, the methodology presented in this section will continue to work even though the change tracking is activated in MS-Office. Enabling the change tracking means that personal information is inserted into the document. Therefore, the Document Inspector signals the presence of the change tracking as an anomaly and proceeds to eliminate this information from the document.

Example 2 (Coding with rsid). As an example, it can be considered that the document under scrutiny has 19 occurrences of the *rsid* values:

```
<w:p w:rsidR="00463DF8" w:rsidRDefault="00463DF8" w:rsidP="00463DF8">
<w:r w:rsidRPr="0074047B">
<w:p w:rsidR="00463DF8" w:rsidRDefault="00463DF8" w:rsidP="00463DF8">
<w:r w:rsidRPr="008C3D74">
<w:r w:rsidRPr="0074047B">
<w:p w:rsidR="00463DF8" w:rsidRPr="008C3D74" w:rsidRDefault="00463DF8" w:rsidP="00463DF8">
<w:p w:rsidR="000E634E" w:rsidRPr="00463DF8" w:rsidRDefault="00463DF8">
<w:sectPr w:rsidR="000E634E" w:rsidRPr="00463DF8" w:rsidSect="009B2A88">
```

Thus, it has 152 (19x8) characters to store information (see Table 4).

Let the message “*this message is hidden in a word document*” (41 characters) to be hidden. The first step is to replace every character of the message with the 2 characters that are the relative representation of the ASCII code (see Table 5).

A sequence of 82 characters is obtained, with a further 70 symbols “0” attached. Thus, a string of 152 symbols is obtained (see Table 6).

Table 4. Sequence of *rsid* values

00 46 3D F8	00 46 3D F8	00 46 3D F8	00 74 04 7B	00 46 3D F8
00 46 3D F8	00 46 3D F8	00 8C 3D 74	00 74 04 7B	00 46 3D F8
00 8C 3D 74	00 46 3D F8	00 46 3D F8	00 0E 63 4E	00 46 3D F8
00 46 3D F8	00 0E 63 4E	00 46 3D F8	00 9B 2A 88	

Table 5. Coded message

<i>t h i s</i>	<i>m e s</i>	<i>s a g e</i>	<i>i s</i>	<i>h i d d</i>
74 68 69 73	20 6D 65 73	73 61 67 65	20 69 73 20	68 69 64 64
<i>e n i n a</i>	<i>w o r d</i>	<i>d o c</i>	<i>u m e n</i>	
65 6E 20 69	6E 20 61 20	77 6F 72 64	20 64 6F 63	75 6D 65 6E
<i>t</i>				
74				

Table 6. Sequence of *rsid* values with hidden data

74 68 69 73	20 6D 65 73	73 61 67 65	20 69 73 20	68 69 64 64
65 6E 20 69	6E 20 61 20	77 6F 72 64	20 64 6F 63	75 6D 65 6E
74 00 00 00	00 00 00 00	00 00 00 00	00 00 00 00	00 00 00 00
00 00 00 00	00 00 00 00	00 00 00 00	00 00 00 00	

Finally it will be enough to replace, in an XML file, the string of symbols in Table 6 to the values of the *rsid* attributes in order to complete the steganography process.

```
<w:p w:rsidR="74686973" w:rsidRDefault="206D6573" w:rsidP="73616765">
<w:r w:rsidRPr="20697320">
<w:p w:rsidR="68696464" w:rsidRDefault="656E2069" w:rsidP="6E206120">
<w:r w:rsidRPr="776F7264">
```

```
<w:r w:rsidRPr="20646F63">
<w:p w:rsidR="756D656E" w:rsidRPr="74000000" w:rsidRDefault="00000000" w:rsidP="00000000">
<w:p w:rsidR="00000000" w:rsidRPr="00000000" w:rsidRDefault="00000000">
<w:sectPr w:rsidR="00000000" w:rsidRPr="00000000" w:rsidSect="00000000">
```

Obviously the message to be hidden would be preferably encrypted before embedding it (see Section 7). □

5 Data Hiding by Zero Dimension Image

The methodology proposed in this section uses an OLE-object (of type “image”), inserted into a MS-Office document in order to contain the information to be hidden. This object, which is totally compatible with the OOXML standard, will:

- be located in the upper-left position and placed in any of the pages that make up the document;
- have both the height and width equal to zero;
- be placed “behind the text”.

These properties will make it possible to hide the image during the visualization or modification of the document. It is worth noting that the file associated to OLE-object, even if declared as “image”, can in reality be any type of file (text, audio, etc.) with a appropriate extension (.jpg, .bmp, etc.). Therefore, this methodology can be used in order to hide data of a different nature, and is not only limited to images. The identification of the OLE-object and the decoding of the hidden text make it more difficult to associate files of reduced dimensions and encrypt the message to be hidden.

A simple and fast method to hide information using this methodology is the following:

- rename the file which contains the hidden message with an extension compatible with an image type;
- insert the image introduced in the previous step into the Word, Excel or PowerPoint document;
- modify the layout of the text related to the image, setting the “Behind the text” style;
- move the image to the upper-left position;
- from the menu “Dimension and position” set both the height and width of the image to 0.

The folder where to copy the OLE-object associated to the file varies according to the type of MS-Office document worked on, with it being *word\media* for Word files, *xl\media* for Excel files, and *ppt\media* for PowerPoint files.

Another way of applying such methodology is to work directly on the XML files. In this case, it is necessary – besides copying the file containing the message to hide in the proper directory (of the ZIP container) – to insert in the XML files the elements to:

- relate to the image;
- declare the presence of the image;
- set the position of the image on the upper-left;
- set the image placed behind the text;
- set the dimensions of the image equal to zero.

In order to set the dimensions of the image to zero, the XML `extent` attribute will have to be worked on (see pp. 3173-3176 in the ECMA-376 specifications [3]). This element, in fact, defines the dimension of the bounding box that contains the image. Therefore, reducing the height and width of the bounding box to zero, will obtain the desired effect. Two examples of the `extent` element, respectively for Word and Excel files, are the following:

```
<wp:extent cx="0" cy="0" />
<a:ext cx="0" cy="0" />
```

Where attributes `cx` and `cy` are, respectively, the width and height of the bounding box. In Excel files, among the elements used to describe the image inserted in the spreadsheet, there are:

```
<xdr:from>
  <xdr:col>0</xdr:col>
  <xdr:col0ff>9525</xdr:col0ff>
  <xdr:row>0</xdr:row>
  <xdr:row0ff>28575</xdr:row0ff>
</xdr:from>
<xdr:to>
  <xdr:col>0</xdr:col>
  <xdr:col0ff>161925</xdr:col0ff>
  <xdr:row>0</xdr:row>
  <xdr:row0ff>28575</xdr:row0ff>
</xdr:to>
```

These elements identify the box of cells that contains the image (see pp. 3516-3517, 3523-3524 and 3532-3533 of the ECMA specifications [3]). The coordinates (line, column) are relative to the two cells situated respectively in the upper-left and lower-right corners. Therefore, in order to reduce the dimensions of the image to zero, it is sufficient to reduce the box of cells that contains it (`<xdr:col>0` and `<xdr:row>0`) to zero. Thus, there is no need to place the image in the upper-left position due to it already being not in a selectable position: the cell with the coordinates (0,0).

In order to set the image in the upper-left corner position of the page, for Word files, it will be necessary to operate on the `position` element (see pp. 3480-3483 of the ECMA specifications in [3]). This element indicates the position of the image with respect to a part of the document (page, column, paragraph). Therefore, placing the image at a distance 0 of the “page” will obtain the desired effect. An example of how the block of elements on which the modification operates, is the following:

```
<wp:positionH relativeFrom="column">
  <wp:posOffset>1685925</wp:posOffset>
  <wp:positionV relativeFrom="page">
  <wp:posOffset>967105</wp:posOffset>
```

The attribute `relativeFrom` indicates the part of the document in relation to which the position will be calculated while `posOffset` is the position. Therefore, upon placing the image on the left, the following elements will be modified:

```
<wp:positionH relativeFrom="page">
<wp:posOffset>0</wp:posOffset>
<wp:positionV relativeFrom="page">
<wp:posOffset>0</wp:posOffset>
```

In order to place the image in the upper-left position, the `<a:off x="0" y="0"/>` element cannot be used due to the position indicated by the x and y coordinates referring to the paragraph and not to the page.

There is a problem for PowerPoint files, where the image, also if reduced to dimension zero and placed in the upper-left corner position, could still be selected by using the “Select Area” function. Moreover, it is not possible to insert an image outside a slide. In fact, the image would be interpreted as an anomaly by the Document Inspector. This methodology, therefore, is not really suitable for PowerPoint files.

6 Data Hiding by Office Macro

A macro is a group of commands which make it possible to obtain a series of operations with a single command [14]. Thus, a macro is a simple recording of sequence of commands which are already available in a software. For this reason, there would seem no need for a programming language. However, macro has acquired a programming language that, in the event of MS-Office, is Visual Basic. The new format of MS-Office, as previously stated, in order to guarantee a greater level of security does not allow macro to be saved inside the file. When using macro in documents, it is necessary to enable this function as well as modify the extension of the name file, which will be: *.docm*, *.xslm*, *.pptm*, etc.. The structure of the files with macro (e.g. *example.docm*) and without (e.g. *example.docx*) is different. This is evident when carrying out a simple test: changing the extension of the file from *.docm* to *.docx* and displaying the document, the system gives an error message indicating that the format is not the one expected. However, MS-Office can open the file, recognizing it as a document with macro and processing it as a normal *.docm* file.

Thus, it is possible to consider using MS-Office macro as a channel to transmit hidden information. In fact, macro can be seen as a function $F(x)$ with $x \in X$, where X is the set of the inputs. Therefore, it is possible to hide information:

- in the description of the function $F(x)$;
- in the value associated to the function $F(k)$, where $k \in K$ and $K \subseteq X$ is the set of stego-key that are highly unusual inputs.

In the first case, the information to be hidden will be stored inside the macro. For example, it is possible to insert the data to be hidden as a comment to the code or as a value assigned to a variable.

In the second case, as consequence of specific input, the macro has a behavior which generates an output that makes the hidden data visible. As an example, in a Word document, a macro given a word as input, searches for it in the text and highlights it in yellow. There could be another routine in the macro, that can only be executed if the searched word is the stego-key, than highlights several characters in the document in yellow. These characters, read in sequence, are the hidden information. In this case:

- the macro will be recognized as reliable by a user as it carries out the task for which it has been realized;
- inside the code, the characters of the hidden message will not be explicitly present but only the coordinates of the corresponding position in the document;
- only who knows the right stego-key will know the secret.

This methodology does not place limits on the amount of information that can be hidden. In fact, a macro does not pre-exist but is created or modified according to the data to be hidden.

7 Methodologies Compared

The Document Inspector, as indicated in Section 2, is the tool supplied by Microsoft to search for and remove any eventual information hidden in MS-Office files. Thus, for an Information Hiding methodology to be considered good, it has to pass the controls of this tool. All four methodologies presented in this paper pass the analysis of the Document Inspector. In addition to controlling and removing hidden information with the Document Inspector, MS-Office also carries out a type of optimization and normalization of the ZIP container every time the file is saved. These operations consist of eliminating everything that it is not recognized as valid for the application (e.g. files attached without a link) as well as reorganizing the elements that reorganize the XML code according to a predetermined layout. As a consequence, the techniques presented in Sections 3 and 4 are vulnerable. In fact, as a result of a save action, MS-Office compresses all files in a ZIP container using the default algorithm (DeflateS) and assigns new values to the rsid attributes. Therefore, in order to avoid that the hidden information is removed as a result of an “involuntary” save action (e.g. automatic saving), it is worthwhile marking the document as the “final version”.

Let us make some considerations about the overhead of the proposed methods. In the method discussed in Section 3, the overhead is a function of the compression ratio of the different algorithms. Therefore, the dimension of the file can either increase, remain unchanged or diminish. On the other hand, the method presented in Section 4 has a null overhead if the text to be hidden is less than $32 \cdot T$. In other words, it is necessary to increase T by adding XML

lines, in order to host longer secret information. The overhead introduced by the solution proposed in Section 5 is a function of two values: the size of the attached image file, that contains the hidden data and the elements added in the XML files by the technique. Finally, in the case discussed in Section 6, the overhead introduced is a function of the size of the macro code. The four methods discussed in this paper can all be applied simultaneously to the same document. In conclusion, in order to obtain confidentiality, the data to be embedded should be encrypted using a symmetric key algorithm.

8 Conclusions

MS-Office documents are the most used file type around the Internet. This makes such file format the ideal candidate for implementing an effective steganographic scheme. This paper proposes four new methods for hiding data in MS-Office documents. The common feature of all the methods is that they pass the Document Inspector analysis. The first technique exploits the chance to use several compression algorithms when creating/managing MS-Office documents. The second technique uses the possibility to conceal information in the values of some special attributes which are normally used to specify unique identifiers when tracking different editing sessions. Using the third method it is possible to hide data by using special sized images embedded in a MS-Office document. The fourth and last technique takes advantages from the adoption of macros for implementing data hiding. In order to increase the information to be concealed in a MS-Office document, all the four methods can be used together.

References

1. Castiglione, A., De Santis, A., Soriente, C.: Taking advantages of a disadvantage: Digital forensics and steganography using document metadata. *Journal of Systems and Software* 80(5), 750–764 (2007)
2. Deutsch, P.: DEFLATE Compressed Data Format Specification version 1.3 (May 1996), <http://www.ietf.org/rfc/rfc1951.txt>
3. ECMA International: Final draft standard ECMA-376 Office Open XML File Formats - Part 1. In: ECMA International Publication (December 2008)
4. Erika Ehrli, M.C.: Building server-side document generation solutions using the open xml object model (August 2007), <http://msdn.microsoft.com/en-us/library/bb735940%28office.12%29.aspx>
5. Frank Rice, M.C.: Microsoft MSDN. Introducing the Office (2007) Open XML File Formats (May 2006), <http://msdn.microsoft.com/it-it/library/aa338205.aspx>
6. Garfinkel, S.L., Migletz, J.J.: New xml-based files implications for forensics. *IEEE Security & Privacy* 7(2), 38–44 (2009)
7. Hao-ran, Z., Liu-sheng, H., Yun, Y., Peng, M.: A new steganography method via combination in powerpoint files. In: 2010 International Conference on Computer Application and System Modeling (ICCA SM), vol. 2, pp. V2-62–V2-66 (October 2010)

8. Jing, M.Q., Yang, W.C., Chen, L.H.: A new steganography method via various animation timing effects in powerpoint files. In: 2009 International Conference on Machine Learning and Cybernetics, vol. 5, pp. 2840–2845 (July 2009)
9. Lin, I.C., Hsu, P.K.: A data hiding scheme on word documents using multiple-base notation system. In: 2010 Sixth International Conference on Intelligent Information Hiding and Multimedia Signal Processing (IIH-MSP), pp. 31–33 (October 2010)
10. Liu, T.Y., Tsai, W.H.: A new steganographic method for data hiding in microsoft word documents by a change tracking technique. *IEEE Transactions on Information Forensics and Security* 2(1), 24–30 (2007)
11. Microsoft Corporation: Compare office professional plus 2010, and the 2007 suite, <http://office.microsoft.com/en-us/professional-plus/professional-plus-version-comparison-FX101871482.aspx> (visited May 2011)
12. Microsoft Corporation: Remove hidden data and personal information from office documents, <http://office.microsoft.com/en-us/excel-help/remove-hidden-data-and-personal-information-from-office-documents-HA010037593.aspx> (visited May 2011)
13. Microsoft Press Release: Microsoft office 2010 now available for consumers worldwide, <http://www.microsoft.com/presspass/press/2010/jun10/06-152010officelaunchpr.mspx> (visited May 2011)
14. MSDN Library: Introduction to macros, <http://msdn.microsoft.com/en-us/library/bb220916.aspx> (visited May 2011)
15. Park, B., Park, J., Lee, S.: Data concealment and detection in microsoft office 2007 files. *Digital Investigation* 5(3-4), 104–114 (2009)
16. Wikipedia: ZIP (file format), http://en.Wikipedia.org/wiki/ZIP_file_format (visited May 2011)