

Embedded Malware Detection using Markovian Statistical Model of Benign Files[★]

M. Zubair Shafiq, Syeda Momina Tabish and Muddassar Farooq

*Next Generation Intelligent Networks Research Center (nexGIN RC)
National University of Computer & Emerging Sciences (NUCES)
Islamabad, Pakistan.*

Syed Ali Khayam

*School of Electrical Engineering & Computer Science (SEECS)
National University of Sciences & Technology (NUST)
Rawalpindi, Pakistan.*

Abstract

Embedded malware is a recently discovered security threat that allows malcode to be hidden inside a benign file. This hidden malcode can be remotely executed using a simple trigger program. It has been shown that embedded malware is not detected by commercial antivirus software even when the malware signature is present in their antivirus database. Therefore, embedded malware can be, more suitably, termed as a ‘disaster in disguise’. In this paper, we present a novel anomaly detection scheme to detect embedded malware. We first analyze byte sequences in benign files to show that benign files’ data generally exhibit a 1-st order dependence structure. Consequently, conditional n -grams provide a more meaningful representation of a file’s statistical properties than traditional n -grams. To capture and leverage this correlation structure for embedded malware detection, we model the conditional distributions as *Markov n -grams*. For embedded malware detection, we use an information-theoretic measure, called entropy rate, to quantify changes in Markov n -gram distributions observed in a file. We show that the entropy rate of Markov n -grams gets significantly perturbed at malcode embedding locations, and therefore can act as a robust feature for embedded malware detection. We evaluate the proposed Markov n -gram detector on a comprehensive malware dataset consisting of more than 37 thousand malware samples and 1,800 benign samples of six well-known filetypes. We show that the Markov n -gram detector provides better detection and false positive rates than the only existing embedded malware detection scheme.

1 Introduction

2 Malware sophistication has evolved considerably during the last decade. In
3 particular, due to emerging financial motivations for attackers, malware trends
4 are now shifting towards stealthy attacks. The challenge faced by stealthy mal-
5 code is to reach vulnerable hosts undetected and then to stay undetected on
6 the hosts. ‘The longer a threat remains undiscovered in the wild, the more
7 opportunity it has to compromise computers before measures can be taken to
8 protect against it. Furthermore, its ability to steal information increases the
9 longer it remains undetected on a compromised computer’ [1]. Code obfusca-
10 tion, (self-)encryption and polymorphism are commonly-used code transfor-
11 mations that are used by stealthy malware to avoid detection.

12 In their seminal work, Stolfo et al. discovered a new type of stealthy threat
13 called *embedded malware* [2]. Under this threat, the attacker embeds the ma-
14 licious code or file inside a benign file on the target host. It was shown that
15 embedded malware cannot be detected by the signature-based antivirus detec-
16 tors even if a malware’s exact signature is present in the detector’s database
17 [2], [3]. In fact, intelligently infected files can even be opened by their respective
18 application software without providing any observable hint of the infection.
19 Intelligent embedding can be further enhanced to allow automatic execution
20 of embedded malcode when the benign file is opened [3]. Embedded malware
21 is potentially a serious security threat and accurate anomaly detection tech-
22 niques must be developed to mitigate it.

23 In this paper, we propose a novel statistical anomaly detection scheme for
24 embedded malware detection. Using correlation analysis, we first show that
25 benign files exhibit a clear 1-st order dependence structure which can be mod-
26 eled using Markov chains. We therefore propose to characterize the statistical
27 properties of a benign file using conditional n -gram distributions, referred to
28 as *Markov n -grams*, instead of the traditional n -grams. For embedded mal-
29 ware detection, we compute running Markov n -grams over non-overlapping

* This paper is an extended version of our earlier work, [7], available at <http://www.springerlink.com/content/978-3-540-70541-3/>

Email addresses:

zubair.shafiq,momina.tabish,muddassar.farooq@nexginrc.org (M. Zubair Shafiq, Syeda Momina Tabish and Muddassar Farooq), khayam@niit.edu.pk (Syed Ali Khayam).

30 windows in a file. We then use an information-theoretic measure, called *en-*
31 *trophy rate*, to quantify perturbations in the Markov n -grams due to embedded
32 malware. The results of our experiments show that the entropy rate of Markov
33 n -grams gets significantly perturbed at malware embedding locations. For au-
34 tomated detection, we observe that the aggregate entropy rate distribution of
35 benign files approaches Gaussianity for large training samples¹. Therefore, a
36 statistical range of benign entropy rates can be defined using the parameters
37 of the baseline Gaussian distribution. Entropy rate values outside this range
38 can then be classified as malicious.

39 We compare the proposed Markov n -gram detector with the only known em-
40 bedded malware detector [2] using two comprehensive and diverse infected
41 datasets. The first dataset is created by randomly embedding malware into
42 benign files. The second dataset is created by randomly embedding naively
43 encrypted malware into benign files. Both datasets are generated from 1,800
44 benign samples (including JPG, MP3, ZIP, EXE, PDF and DOC files) and 37,420
45 malware samples (containing viruses, worms trojans, spyware, and exploit
46 codes). We show that the Markov n -gram detector consistently outperforms
47 the only existing embedded malware detector [2] in terms of both detection
48 and false positive rates. In comparison to commercial-of-the-shelf (COTS) an-
49 tivivirus (AV) software, our detector provides a significantly higher detection
50 rate at the cost of higher false positive rates. Therefore, we argue that, due to
51 their complementary strengths, very high accuracy can be achieved when the
52 Markov n -gram detector is deployed in conjunction with COTS AV software.
53 In this study, we also explore this proposition.

54 **Organization of the Paper.** The rest of the paper is organized as follows.
55 We discuss in detail the infected datasets created for this study in Section 2.
56 We present realistic attack scenarios, using different filetypes, for embedded
57 malware in Section 3. In Section 4, we propose our Markov n -gram detector
58 and in Section 5 we compare the detection accuracy of our proposed detector
59 with other relevant techniques and state-of-the-art antivirus products using
60 the infected datasets. In Section 6, we discuss the limitations of the proposed
61 Markov n -gram detector. We provide an overview of related research in the
62 field of embedded malware detection and n -gram based techniques in Section
63 7. Finally, we conclude the paper with an outlook to our future work.

64 2 Dataset

65 In this section, we first describe the benign and malware datasets used in this
66 paper. We then introduce our tool **NERGAL** that embeds any given infection at

¹ This is a direct consequence of the central limit theorem.

Table 1
 Statistics of Benign Files used in this Study

Filetype	Quantity	Average Size (kilo-bytes)	Minimum Size (kilo-bytes)	Maximum Size (kilo-bytes)
JPG	300	1,097.8	3	1,629
MP3	300	3,384.4	654	6,210
ZIP	300	1,489.6	8	9,860
EXE	300	4,095.0	48	15,005
PDF	300	1,513.1	25	20,188
DOC	300	1,015.2	44	7,706

any random location within a benign file. Using this tool, we produce a large embedded malware dataset².

2.1 Benign Dataset

The benign dataset for our experiments consists of six different filetypes: JPG, MP3, ZIP, EXE, PDF and DOC. These filetypes encompass a broad spectrum of commonly used files ranging from compressed to redundant and from executables to document files. Each set of benign files contained 300 typical samples of the corresponding filetype, which provide us a total of 1,800 benign files. We ensured the generality of the benign dataset by randomizing the sample sources. More specifically, we queried well known search engines with random keywords to collect these files. In addition, we also collected typical samples on the local network of our virology lab.

Some pertinent statistics of the benign dataset used in this study are tabulated in Table 1. It can be observed from Table 1 that the benign files have diverse sizes varying from 3 KB to 20 MB, with an average file size of approximately 2 MB. We show later in the paper that this diversity in file sizes provides valuable insights into an important aspect of embedded malware detection, that is, whether or not a detector is able to detect the embedded malware in large files where the statistical contents of the malicious code are simply averaged out.

The executable files collected for this study include both compiled and compressed (installation) executables. The JPG, MP3, and ZIP file formats are inherently compressed so the n -grams on the data portion of these files should provide distributions that are fairly uniform. Evaluation and detection of embedded malware in these uniform distributions is an important issue which was originally raised in [2].

² Some characteristic samples of the dataset and the tool, NERGal, are publicly available at <http://www.nexginrc.org>.

Table 2
 Statistics of Malware used in this Study

Major Category	Minor Category	Quantity	Average Size (kilo-bytes)	Minimum Size (bytes)	Maximum Size (kilo-bytes)
Backdoor	Win32	3,444	285.6	56	9,502
Constructor	DOS	178	104.2	62	7,241
Constructor	Win32	172	398.5	371	5,971
Email Flooder	-	148	343.5	1,430	4,262
Email Worm	Win32	935	73.5	148	762
Exploit	-	242	101.1	370	1,912
Flooder	-	154	168.1	486	981
IRC Worm	-	485	34.4	56	1,072
Nuker	-	140	188.1	4,000	680
Trojan	BAT	649	20.2	12	708
Trojan	DOS	971	27.0	4	1,818
Trojan	Win32	983	125.4	12	2,998
Virus	Boot	1,514	32.5	108	1,490
Virus	DOS	16,236	18.7	5	1,860
Virus	MS Office	2,596	53.5	118	4,980
Virus	Win32	991	44.3	175	1,018
Worm	Win32	153	110.5	97	2,733

93 *2.2 Malware Dataset*

94 Malware samples, especially recent ones, are not easily available on the Inter-
 95 net. Computer security corporations do have an extensive malware collection,
 96 but unfortunately they do not share their malware databases on the Internet.
 97 We could only locate ‘VX Heavens Virus Collection’ [14] which is available for
 98 free download in the public domain. This is a comprehensive database that
 99 contains a total of 37,420 malware samples. The sample consists of backdoors,
 100 constructors, flooders, nukers, sniffers, droppers, spyware, viruses, worms and
 101 trojans etc.

102 A detailed description of the malware used in our study is provided in Table
 103 2. The average malware size in this dataset is 64.2 KB. Note that this size is
 104 significantly larger than the average size (2 MB) of the benign files. Moreover,
 105 the sizes of malware samples used in our study vary from 4 bytes to more than
 106 14 MB. Clearly, small sized malware are harder to detect than larger ones.

107 *2.3 Infected Dataset*

108 We developed an inhouse software tool, called NERGAL, that could insert an
 109 infection into benign files at any given location in the benign file. The algo-

Table 3
 NERGal Algorithm

input: benign file, malicious payload
 output: infected file with embedded malware

Read benign file in a byte array C
 Read malicious payload in a byte array M
 Create a infected-data byte array S such that $|S| = |C| + |M|$, where $|X|$ is the size of array X
 Copy C 's header to S : $S[1 \dots h] \leftarrow C[1 \dots h]$, where h is the size of C 's header
 Generate a random number r : $h + 1 \leq r \leq |C|$
 $S[h + 1 \dots r] \leftarrow C[h + 1 \dots r]$
 $S[r + 1 \dots r + |M|] \leftarrow M[1 \dots |M|]$
 $S[r + |M| + 1 \dots |C| + |M|] \leftarrow C[r + 1 \dots |C|]$

Write S to the infected file

110 rithm of NERGal is shown in Table 3. NERGal ensures that the infections are
 111 inserted after the header of the benign files to avoid filetype corruption. The
 112 tool also generates a detailed infection report, which provides details about
 113 the sample malware that was used to infect each benign sample and its offset
 114 in each sample.

115 We have created two infected datasets for this study. The first infected dataset
 116 is created by simply embedding malware inside benign files. The second in-
 117 fected dataset is created by encrypting the malware before embedding. We
 118 use the ROT-13 Caesar cipher for malware encryption. While more sophis-
 119 ticated encryption techniques are certainly possible, we use a simple substi-
 120 tution cipher because it does not alter the inherent statistical properties of
 121 the malcode. Therefore, while COTS AV software will not be able to detect
 122 this naively encrypted malcode, we can intuitively argue that the accuracy of
 123 statistical anomaly detection techniques should remain unaffected under this
 124 simple encryption. (We show later that this is not the case.)

125 The complete virus dataset is used for every filetype mentioned in the benign
 126 dataset. Therefore, the embedded malware dataset for each filetype consists of
 127 37,420 files and the total number of files in both infected datasets are 449,040.
 128 The average file size in both datasets is 2,267.5 KB.

129 3 Attack Scenarios

130 In this section, we will discuss potential real-world attack scenarios that can
 131 be realized using embedded malware. The potential threat of the embedded
 132 malware depends on its stealthiness. There are two possible ways in which its
 133 stealthiness can be compromised: 1) detection by anti-virus software, and 2)
 134 error message on opening or degraded quality. As explained by the authors in
 135 [3], and independently verified by us, embedded malware cannot be detected

Table 4
Results for opening

Cannot Open without Error Message	Cannot Open with Error Message	Open with Error Message	Open with Degraded Quality	Open without Error Message
JPG				
10.0%	0.0%	0.0%	88.3%	1.6%
MP3				
0.0%	0.0%	0.0%	0.0%	100.0%
ZIP				
1.1%	3.7%	47.9%	47.4%	0.0%
EXE				
7.0%	80.0%	1.0%	N/A	12.0%
PDF				
2.4%	4.9%	74.1%	6.5%	12.2%
DOC				
0.0%	89.5%	9.6%	0.0%	0.9%

136 by COTS AV software. Similarly, ‘intelligently’ embedded malware inside be-
 137 nign (document, media or application) files does not affect their integrity as
 138 these infected files continue to open by their respective application software
 139 without any error message or degradation in quality. To this end, we manually
 140 opened the 6,000 infected files and observed how many of them open correctly
 141 after naive (i.e., completely random) malware embedding. The infected files
 142 followed various trends upon opening; some files opened with or without error
 143 messages, other did not open at all, and still some other opened with degraded
 144 quality. The opening of the infected files revealed interesting insights for dif-
 145 ferent filetypes whose stats are summarized in Table 4³. We will now discuss
 146 these insights for every filetype separately.

147 3.1 *JPG files*

148 Most of the **JPG** files are previewed after being infected as shown in Figures
 149 1(b), 1(d) and 1(d); only 10% files are not opened. Most of the infected files
 150 have not incurred any distortion as shown Figure 1(b) and (d). In some cases,
 151 such as Figure 1(f), some distortion is visible in the infected files. Here we
 152 highlight that this image distortion is introduced because we are randomly
 153 embedding malware inside the **JPG** files without any regard to the inherent
 154 structure of the benign file’s content. However, sophisticated watermarking
 155 techniques can be used to minimize the distortion due to malware embedding.
 156 Furthermore, we highlight that an automated detector that does not under-
 157 stand the file’s semantics and the original file will not be able to understand

³ The category of ‘degraded quality’ is not applicable for **EXE** files.



Fig. 1. JPEG files before and after malware embedding.

158 or quantify this distortion in the image quality.

159 3.2 MP3 files

160 In our pilot studies, we have observed that MP3 song files can serve as very
 161 potent carriers of embedded malware; 100% infected MP3 files (with embedded
 162 malware) play from start to finish *without any error messages or perceptual*
 163 *distortion in audio quality*. This came as a surprise to us but we found out
 164 that it is due to the frame-based structure of MP3 files. Each frame in the MP3

165 file format is preceded by a re-sync marker. Corrupt frames (without re-sync
166 markers) are simply bypassed by the media players during playback. Since
167 most Internet song sharing portals use the MP3 file format, an attacker can
168 use simple embedding to infect a benign MP3 file by a malware and then can
169 distribute the infected file via Internet song sharing portals or peer-to-peer file
170 sharing software.

171 3.3 ZIP files

172 The ZIP files are stored in an object-wise manner. Every file in the overall ZIP
173 file has a separate header. This protects the complete ZIP file from getting
174 corrupted. 95% of the ZIP files open with an error message because of failed
175 CRC checksum. However, they extract successfully without any further errors.

176 3.4 EXE files

177 The Microsoft Portable Executable (PE) file has several well-known vulnera-
178 bilities [18]. A random code can be embedded in free spaces between different
179 sections or in the padding portions of the file without violating the integrity
180 of executable file. The author in [18] has also described a number of intel-
181 ligent embedding techniques which can automatically execute the embedded
182 code once the infected file is executed. Our results show that with random
183 embedding, 12% of the files can open without any error message. Since the
184 embedded files are undetected by COTS AV software, an attacker can embed
185 malware inside common benign files – for instance, WINWORD.EXE – and the
186 infected file will go unnoticed through the COTS AV software deployed inside
187 the network or on the host. Such infected file can be transported to different
188 hosts using well-known peer-to-peer file sharing software or by making the file
189 freely available for download. Later on, a user can be tricked into starting a
190 trigger program (in the form of a plug-in or a macro) to launch the malicious
191 code. Examples of similar attacks have recently been reported in [3]–[6]. It
192 should be noted that disabling macros and plug-ins is not a viable option be-
193 cause there are many useful benign programs (e.g., MathType, Adobe PDF
194 printer, flash player, etc.) that are launched as macros or plug-ins.

195 3.5 PDF files

196 The PDF files have an object-oriented encoding format. The random infec-
197 tion only corrupts a particular object and other objects remain intact. Adobe
198 reader notices the corruption due to the infection but the files still open. Most

199 of the PDF files open with or without error message. The error messages typi-
200 cally appear when we scroll down to the document page where the infection is
201 introduced. PDF is a modern document format with ‘object oriented dynamic
202 composability’ which allows the user to include embedded objects such as
203 video clips, wave sounds or bitmap images inside a document. The embedded
204 objects can be invoked by simply clicking on the object. An attacker can create
205 a fake embedded object which, in addition to some benign looking activity,
206 executes the malcode [3].

207 3.6 DOC files

208 The DOC file format also supports embedded objects. However, 90% of the
209 files displayed an error message on double clicking and could not be opened.
210 The authors in [3] have shown that intelligent embedding in DOC files can add
211 embedded objects to a DOC file that can be executed from the DOC file editor.
212 Similar to PDF, DOC is also a modern document format which allows the user
213 to include embedded objects which can be invoked by deceiving a naive user.

214 4 Modeling and Quantification of n -gram Information

215 In this section, at first, we present pilot experimental studies carried out to
216 evaluate previously proposed schemes for embedded malware detection. Then
217 we analyze byte-level correlation structure present in files which leads us to a
218 statistical model (Markov n -gram) of benign byte sequences. We quantify per-
219 turbations in this statistical model using a well-known information-theoretic
220 measure (entropy rate). Entropy rate thresholds are selected for different file-
221 types using the Gaussianity properties of *the sum of sampled entropy rate*
222 *distributions* of benign files.

223 It should be noted that two desirable accuracy objectives of an embedded
224 malware detector are: 1) to detect infected files and 2) to identify the likely
225 location of the embedded infection. We refer to these two objectives as *detec-*
226 *tion* and *location identification*, respectively.

227 4.1 Mahanalobis n -gram Detector

228 In pilot studies, we build on the Mahanalobis n -gram detector, the only scheme
229 designed specifically to detect embedded malware, proposed by the authors in
230 [2]. In [2], the authors proposed to use n -gram analysis for embedded malware

231 detection. An n -gram of a sequence is a normalized frequency histogram (or
232 the distribution) of n bit symbols in the sequence. Stolfo et al. [2] used 1-
233 Centroid, Multi-Centroids and Exemplar files as centroids for modeling benign
234 and malware files. 1-gram and 2-gram distributions were used for this purpose.
235 Mahalanobis distances of a given (unknown) file from the benign and the
236 malware model were used for classification.

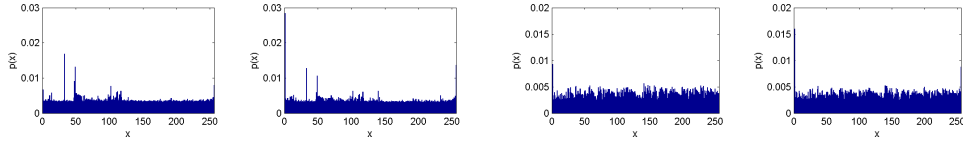
237 One major assumption of the prior study was that the infection appears only
238 at the start or the end of the benign file [2]. Therefore, n -gram analysis was
239 applied only on the truncated files [2]. We argue that this assumption is unre-
240 alistic because it is not capable of detecting embedded malware in the middle
241 of the file. In fact, in our experiments we observed that malware embedded at
242 the start of benign files is detected more frequently by COTS AV software than
243 malware embedded in the middle. Therefore, a pragmatic embedded malware
244 detector should look at the statistical contents of an entire file rather than
245 focussing on a specific location⁴. Mahalanobis n -gram detector provides rea-
246 sonable detection accuracy when the infection appears at the start or the end
247 of a file. However, it cannot provide location identification. We hence revoke
248 the assumption of file truncation and compute n -grams on whole files.

249 Figures 2(a) and 2(b) show the comparison of whole file 1-grams of sample be-
250 nign and infected PDF and JPG files, respectively⁵. It can be clearly seen that
251 no discernable change in the 1-grams is evident in Figures 2(a) and 2(b). It
252 can be intuitively argued that whole file 1-grams of infected files do not change
253 when the size of the benign file is significantly larger than the malware size
254 because the statistical contents of the embedded malware are averaged out by
255 large amounts of benign data. (Recall that the average sizes of the benign and
256 malware files in Tables 1 and 2 were 2 MB and 64.2 KB, respectively.) This sit-
257 uation is quite common because malware are generally designed to have small
258 sizes to make them fit inside buffer overflows/file pads/email attachments or
259 to avoid network-based detection during an initial downloading stage.

260 The authors in [2] also carried out n -gram analysis of a file in a block-wise
261 manner in order to detect the exact location of the embedded malware. Ex-
262 periments were repeated using block sizes of 500 bytes and 1000 bytes. The
263 significance of the block size is that it sets an approximate bound on the
264 minimum size of malware that can be possibly detected. We repeated these
265 block-wise n -gram experiments on our datasets as well. Figure 3 shows some
266 representative results of the Mahalanobis distance between the block-wise 1-

⁴ Here we acknowledge the complexity incurred by n -gram analysis of whole files. Nevertheless, we tradeoff complexity for accuracy throughout this paper. In other words, we expect that the proposed detector will be complemented by signature-based detector.

⁵ Byte value 0 has the highest frequency because of zero padding that is used for block alignment.



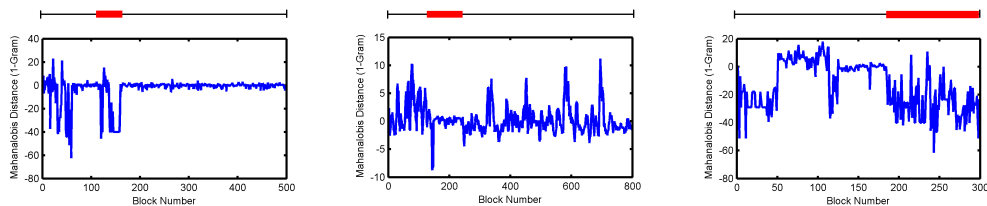
(a) 1-gram histograms of benign and infected PDF file (b) 1-gram histograms of benign and infected JPG file

Fig. 2. Comparison of 1-gram histograms of benign and infected files.

267 gram distribution and the benign file model. We use a block size of 1000 bytes
 268 and plot the Mahalanobis distance between every block and the benign file
 269 model; qualitatively similar results were obtained for other block sizes. One
 270 can see in Figure 3 that the block-wise 1-gram Mahalanobis distance does not
 271 provide significant perturbations that could help in detecting the embedded
 272 malware.

273 As a logical improvement of 1-gram analysis, we repeated our experiments to
 274 analyze the behavior of block-wise Mahalanobis distances of 2-gram distrib-
 275 utions; block size was 1000 bytes. These experiments revealed that, despite
 276 the increased computational complexity, the performance of the Mahanalo-
 277 bis distance based detector did not improve significantly. Here, we can intu-
 278 itively argue that the block size of 1000 bytes does not provide enough data
 279 to compute an effective statistical distribution. Specifically, in case of 2-gram
 280 distribution we only have 1000 data values to fill 65, 536 bins. The ratio (data
 281 values to distribution bins) of about 1 : 65 for 2-gram is in stark contrast to
 282 the ratio of about 4 : 1 for 1-gram using the block size of 1000 bytes. A sim-
 283 ple solution to this problem is to increase the block size. However, as stated
 284 previously, the block size roughly defines the lower bound on the minimum
 285 size of malware that can be detected. Therefore, there is an inherent tradeoff
 286 between the block size and the minimum malware size that can be detected:
 287 increasing the block size means higher false negative rates thereby degrading
 288 the accuracy of the detector. This reason stopped us from increasing the value
 289 of n and we did not extend our study beyond 2-grams.

290 We also tried the *entropy* measure instead of Mahalanobis distance, but it also
 291 failed to accurately detect embedded malware [12], [13]. From the pilot stud-
 292 ies, we can conclude that the Mahalanobis distance or entropy of simple 1- and
 293 2-gram distributions cannot accurately detect embedded malware. This shows
 294 that a simple n -gram distribution does not provide sufficient information to
 295 detect embedded malware. To rectify this shortcoming in the n -gram distrib-
 296 utions, in the following we provide a different method of computing n -grams
 297 in the following sections.



(a) Block-wise 1-gram Mahalanobis distance for an infected EXE file (b) Block-wise 1-gram Mahalanobis distance for an infected PDF file (c) Block-wise 1-gram Mahalanobis distance for an infected DOC file

Fig. 3. Block-wise 1-gram Mahalanobis distance is unable to show significant perturbations in the infected regions. The horizontal thick bars show the location of the embedded malware.

298 4.2 Correlation Structure in File Data

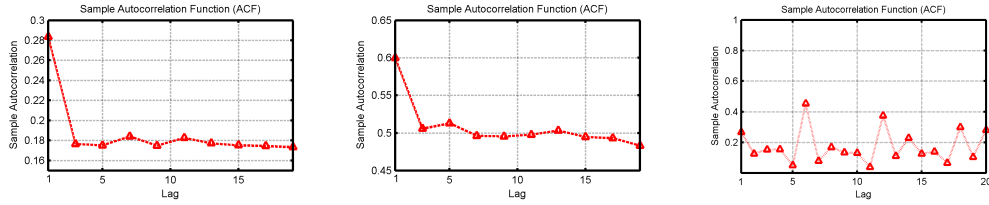
299 We first note that the 2-gram distribution is in fact the joint distribution
 300 of two 1-gram symbols. This joint distribution may contain some redundant
 301 information which is not pertinent to the present embedded malware detec-
 302 tion problem. For accurate detection, it is important that this redundancy
 303 is removed. To this end, we analyzed a number of statistical properties of
 304 the benign files' n -grams. One relevant property that provided us interesting
 305 insights into statistical properties of file data was the analysis of byte level
 306 autocorrelation of benign files.

307 Autocorrelation describes the correlation between the random variables in a
 308 stochastic process at different points in time. For a given lag k , the autocorrela-
 309 tion function of a stochastic process, X_i (where i is the time index) is defined
 310 as:

$$311 \quad \rho[k] = \frac{E\{X_0 X_k\} - E\{X_0\}E\{X_k\}}{\sigma_{X_0} \sigma_{X_k}}, \quad (1)$$

312 where $E\{\cdot\}$ represents the expectation operation and σ_{X_i} is the standard de-
 313 viation of the random variable at time lag i . The value of the autocorrelation
 314 function lies in the range $[-1, 1]$, where $\rho[k] = 1$ means perfect correlation at
 315 lag k (which is obviously true for $k = 0$) and $\rho[k] = 0$ means no correlation at
 316 all at lag k .

317 To observe the level of spatial dependence in the byte sequences of benign files,
 318 we computed their sample autocorrelation functions. Figures 4(a) and 4(b)
 319 show the autocorrelation function plotted versus the lag for EXE and DOC files,
 320 respectively. These autocorrelation results clearly show that the byte sequences
 321 in benign files have 1-st order dependence because the autocorrelation value
 322 takes a fairly significant dip at $k = 2$ and remains constant for higher values
 323 of lag. In other words, once a byte S_i appears, it is more likely that it will be



(a) Autocorrelation results for benign EXE file (b) Autocorrelation results for benign DOC file (c) Autocorrelation results for Code Red II worm file

Fig. 4. Autocorrelation function of byte distributions of benign files shows 1-st order dependence. Autocorrelation function of the byte distribution for Code Red II worm shows that the structure of the 1-st order spatial dependence is disturbed.

324 followed by S_i at the next byte location. Clearly, if we are in a zero padded
 325 region of a benign file, a zero valued symbol is highly likely to be followed by
 326 another zero valued symbol.

327 This 1-st order spatial dependence of benign files has direct implications on the
 328 present embedded malware detection problem mainly because this structure is
 329 not observed in malware files, see Figure 4(c). In fact, instead of the 1-st order
 330 dependence, we can instead observe high correlation at $k = 6, 12,$ and 18 . This
 331 lack of 1-st order spatial dependence of a malware can be easily observed by
 332 examining the signature of the Code Red II Worm given below [11]:

```

333 GET /default.ida?XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
334 XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
335 XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
336 XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
337 XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX%u9090%u6858%ucbd3
338 %u7801%u9090%u6858%ucbd3%u7801%u9090%u6858%ucbd3%u7801
339 %u9090%u9090%u8190%u00c3%u0003%u8b00%u531b%u53ff
340 %u0078%u0000%u00=a HTTP/1.0
  
```

341

342 We can see in the signature of Code Red II that it consists of sub-blocks of
 343 6 bytes, as a result, the high correlation values are observed at $k = 6$ or its
 344 integral multiples.

345 In addition to the CodeRed example shown in Figure 4(c), we also conducted
 346 correlation experiments on other malware and benign filetypes. These corre-
 347 lation results were consistent with the already presented results. We hence
 348 deduce that the 1-st order dependence structure due to zero pads in benign
 349 files is not present in malcode. This result is also intuitive because a main
 350 objective of effective malcode development is to limit the size of the malcode.
 351 (Small sized malware can fit into buffer overflows and can avoid arousing sus-
 352 picion during transmission over the network.) This objective is clearly defeated

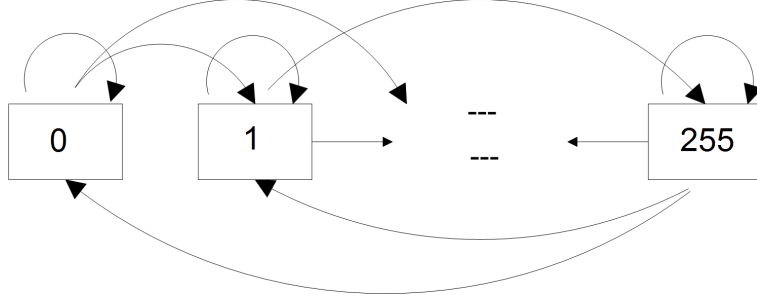


Fig. 5. 1-st order, 256 state, Markov Chain.

353 if an attacker allows large zero pads inside the malcode file.

354 We note that the difference in 1-st order correlation structure of benign and
 355 malicious files is actually a distinguishing feature that can be used to detect
 356 embedded malware. Therefore, in the following section we model and quantify
 357 this distinguishing feature.

358 4.3 A Statistical Model of Benign Byte Sequences

359 We now focus on developing a model for the correlation structure observed in
 360 benign files. Since the correlation shows 1-st order dependence, the underlying
 361 random process (i.e., the byte sequence of benign files in the present context)
 362 can be modeled using an order-1, discrete time Markov chain [13]. Here we note
 363 that a Markov chain characterizes a process in terms of conditional distribution
 364 of its states. For a byte level distribution, a Markov representation simply
 365 implies $2^8 = 256$ conditional probability distributions, each corresponding to
 366 a different byte value. These conditional distributions reduce the size of the
 367 underlying sample space which in the present problem corresponds to removing
 368 redundant information from the joint distribution.

369 The Markov Chain used to model the conditional byte distribution is an order-
 370 1 (256 state) Markov chain as shown in Figure 5. The transition probabilities
 371 are computed by counting the number of times byte i is followed by byte j .
 372 These probabilities can also be expressed as a transition probability matrix. If
 373 the probability of moving from state i to j is $p_{i,j}$, then the transition matrix
 374 for the present problem is given by:

$$375 \mathbf{P} = \begin{bmatrix} p_{0,0} & p_{0,1} & \cdots & p_{0,255} \\ p_{1,0} & p_{1,1} & \cdots & p_{1,255} \\ \vdots & \vdots & \ddots & \vdots \\ p_{255,0} & p_{255,1} & \cdots & p_{255,255} \end{bmatrix}$$

376 Each row of this transition probability matrix provides the conditional dis-
 377 tribution for a distinct byte value. Thus the total number of variables that
 378 characterize this random process (65,536 floating point values) is the same
 379 as the 2-gram distribution. However, these Markov chains provide an alterna-
 380 tive, non redundant and conditional representation of the jointly distributed
 381 2-gram values. Henceforth, we refer to this representation as *Markov n-grams*.

382 We now need an accurate measure that can quantify changes in the Markov
 383 transition probabilities. This measure is presented in the following section.

384 4.4 Quantification of Perturbations in Markov n-grams

385 We need a mathematical measure to quantify changes or perturbations in
 386 the Markov n -gram's transition probability matrix. To this end, we use an
 387 information theoretic measure, called *entropy rate*, which quantifies the time
 388 density of the average information in a stochastic process [13]. Entropy rate
 389 for a sequence of discrete finite random variables X_1, X_2, \dots, X_n is defined as:

$$390 \quad R = \lim_{N \rightarrow \infty} \frac{H(X_1, X_2, \dots, X_n)}{N}, \quad (2)$$

391 where $H(X_1, X_2, \dots, X_n)$ is the joint entropy of random variables X_1, X_2, \dots, X_n .
 392 R does not exist in general. However, for the present n -gram Markov chain
 393 with 256 states, the entropy rate R can be defined in terms of entropy. Let X
 394 be a discrete random variable such that $X = \{x_i, i \in \Delta_n\}$, where Δ_n is the
 395 image of the random variable. Then entropy of X is defined as:

$$396 \quad H(X) = - \sum_{i \in \Delta_n} p(x_i) \log_2 p(x_i). \quad (3)$$

397 Further,

$$398 \quad R = \sum_{i=0}^{255} \pi_i H(X_i), \quad (4)$$

399 where π_i represents the equilibrium probability of being in state i and $H(X_i)$
 400 is the entropy of the conditional distribution of state i (i.e., the entropy of row
 401 i of the transition probability matrix).

402 Asymptotic properties of the entropy rate measure are applicable only in case
 403 of stationary Markov chains [13]. We acknowledge that in general stationarity
 404 will not hold for the present problem. However, the entropy rate expression
 405 does provide us with the *expected entropy* of a discrete time Markov chain.

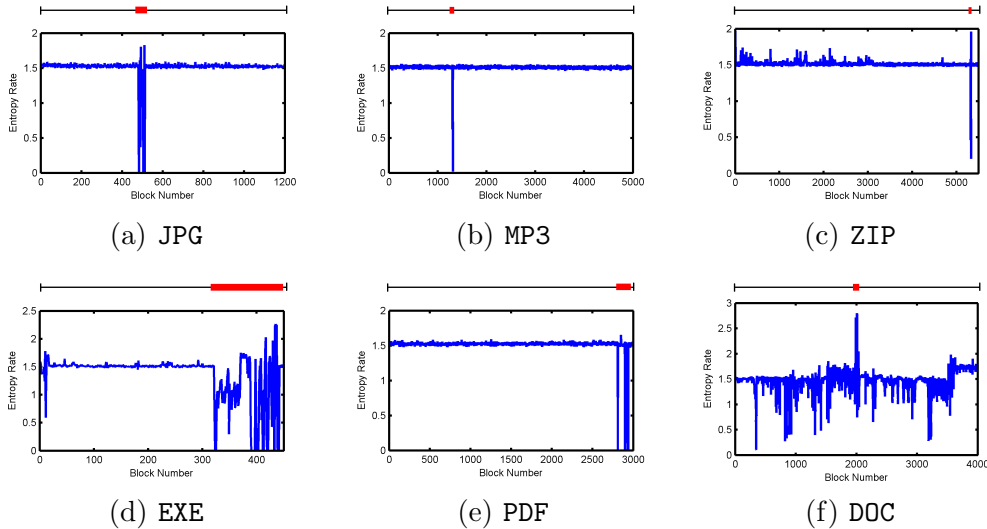


Fig. 6. Entropy Rate of infected files. The horizontal thick bars show the location of the embedded malware.

406 Since we rely on the premise that the statistical properties of the embedded
 407 malware will be different from the statistical properties of the benign file in
 408 which it is embedded, *expected entropy* of the consequent Markov chain (de-
 409 rived from the infected file) should be perturbed at the embedding locations.

410 Figure 6 shows the entropy rate of infected files of every filetype used in
 411 our study. It is clear from Figure 6 that the perturbation is more profound
 412 as compared to those obtained using 1-gram Mahanalobis distance or 2-gram
 413 Mahanalobis distance. This clearly verifies our earlier hypothesis that the con-
 414 ditional distribution discards the redundant information contained in the joint
 415 n -gram distribution, thus providing us a compact representation of the file
 416 data.

417 The results of Figure 6 show that the entropy rate of Markov n -grams can
 418 quantify and highlight perturbations at the locations of the embedded mal-
 419 ware. Thus this measure satisfies the detection and location identification ob-
 420 jectives that we have set for an effective embedded malware detector. However,
 421 for automated detection, we must threshold entropy rate values above and be-
 422 low which an infection would be detected. The following section provides a
 423 flexible yet accurate method of defining this threshold.

424 4.5 Classification using Entropy Rate Thresholding

425 For classification purposes, we need to set an appropriate threshold value on
 426 the block-wise entropy rate values. For this purpose, we develop a generic
 427 model of block-wise entropy rate values in the benign files. During our pilot

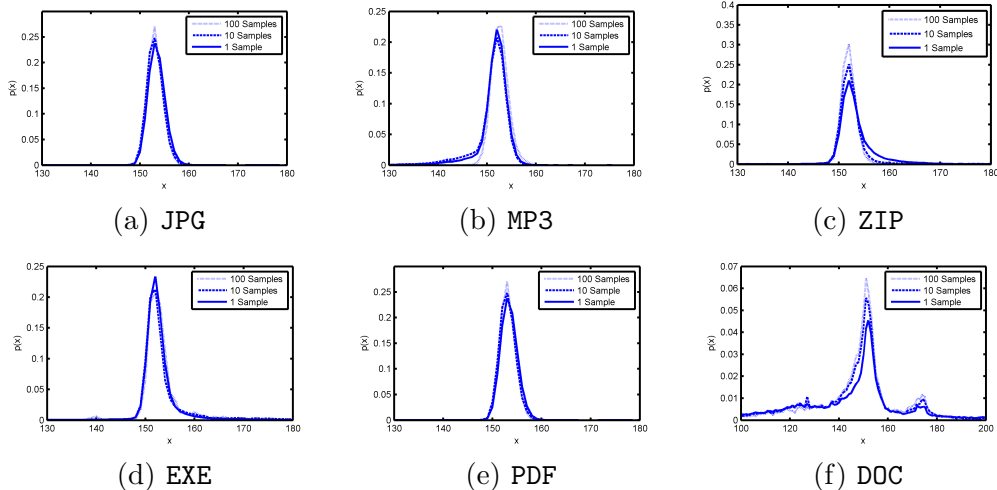


Fig. 7. Sampled entropy rate distributions for different filetypes.

428 studies, we observed that the block-wise values of entropy rates varied in the
 429 range of $[0, 3]$. Therefore, we generated an entropy rate histogram using 300
 430 equal sized bins. We normalized this histogram to obtain the *sampled entropy*
 431 *rate distribution*. Figure 7 shows that the sum of sampled entropy rate distri-
 432 butions approaches Gaussianity as the number of samples (i.e., the benign
 433 files used for training) approaches infinity. This is a consequence of the central
 434 limit theorem which asserts that, for independent finite variance entropy rate
 435 distributions, an aggregated distribution should be normally distributed.

436 Normal distribution is completely specified by its first and second central
 437 moments: mean (μ) and variance (σ^2). Since 99.99% of the times a normal
 438 distribution does not deviate from its mean by more than 5 standard devia-
 439 tions, we can set the upper and lower detection thresholds as: $\eta_{low} = \mu - 5\sigma$
 440 and $\eta_{high} = \mu + 5\sigma$, respectively. Moreover, we integrate the sampled entropy
 441 rate distribution of a test file outside these points to obtain area at the fringes
 442 of the distribution and set a classification threshold k on this area. If the area
 443 outside the $[\eta_{low}, \eta_{high}]$ range is greater than k then the test file is classified as
 444 malicious. Conversely, if the area inside the $[\eta_{low}, \eta_{high}]$ range is less than or
 445 equal to $(1 - k)$ only then an alarm is raised. The value of this threshold was
 446 tuned for the best performance in the ROC space using a randomly sampled
 447 training dataset which was 5% of the total testing dataset [19]. Suitable values
 448 of this threshold were different for different filetypes. The highest value of k
 449 was observed for the DOC filetype. An intuitive feel for the high value of k for
 450 DOC files can be developed with the help of Figure 7. DOC filetype shows worst
 451 convergence to Gaussianity, i.e., significant amount of area is present at the
 452 fringes (see Figure 7(f)). This logically leads us to set a relatively higher value
 453 of the threshold (k) for DOC files.

454 Based on the results of this section, we conclude that entropy rate of Markov
 455 n -grams can achieve both accuracy objectives (i.e., detection and location

456 identification) that we expect from an embedded malware detector. To the
457 best of the authors' knowledge, the location identification objective cannot be
458 achieved by any existing embedded malware detector. Moreover, and again
459 in contrast to existing techniques [2], [3], malware data is *not* required to
460 train our proposed detector. This makes the Markov n -gram detector a true
461 anomaly detector that detects maliciousness by flagging deviations from a
462 robust model of normal behavior. In addition to these desirable properties, the
463 following section shows that the proposed detector provides better accuracy
464 than existing schemes.

465 5 Classification Results

466 As mentioned in Section 2, we perform classification on two infected datasets,
467 each consisting of 224,520 infected files created by infecting benign files of six
468 common types: JPG, MP3, ZIP, EXE, PDF, DOC. For the training of our proposed
469 scheme, we use 5% of the benign dataset.

470 Table 5 provides the detection rates of 3 fully updated⁶ commercial antivirus
471 products: McAfee Antivirus [15], AVG Antivirus [16] and Kaspersky Antivirus
472 [17], Mahalanobis n -gram detector by the authors in [2] and our Markov
473 n -gram detector⁷. We emphasize that embedded malware's signatures are
474 present in the COTS antivirus' databases i.e., the malware are detected as
475 standalone files or exploit codes by the COTS antivirus software. The results
476 tabulated in Table 5 reaffirm that: *commercial antivirus products are not effective in detecting embedded malware*. Moreover, and as expected, the detection
477 rate for all COTS AV software degrade to 0% for the encrypted dataset. Thus,
478 after simple substitution-type encryption, COTS antivirus completely fail to
479 detect the embedded malware in the infected files, irrespective of the type of
480 file being infected. In the un-encrypted infected dataset, some of the embedded
481 malware are detected by COTS antiviruses. An important observation from
482 the experiments is that *the malware detected by COTS antiviruses are mostly embedded at the start of the infected files*. It can also be observed from the
483 Table 5 that code embedding inside EXE and PDF files is detected more often
484 by COTS than other filetypes. While the detection algorithms of these COTS
485 antivirus software are proprietary, we can deduce some properties of these al-
486 gorithms based on our experimental results. For instance, it is obvious that
487 EXE files are scrutinized/scanned more thoroughly than other filetypes. This
488 must be because EXE files are stand alone, execution-ready files which can be
489
490

⁶ by January 2008.

⁷ We also wanted to compare our proposed scheme with the static detection approach proposed of [3]. However, it was not possible because their approach is specific to Microsoft Word and similar document formats.

Table 5
 Detection (TP) rate and False Positive (FP) rate of Antivirus Software and Anomaly Detectors

	McAfee Antivirus [15]	AVG Antivirus [16]	Kaspersky Antivirus [17]	Mahalanobis <i>n</i>-gram Detector	Markov <i>n</i>-gram Detector	Percentage Improvement
<i>unencrypted</i> JPG						
TP rate	0.0%	0.0%	0.0%	76.3%	95.4%	19.1%
FP rate	0.0%	0.0%	0.0%	35.7%	2.7%	33.0%
<i>encrypted</i> JPG						
TP rate	0.0%	0.0%	0.0%	68.9%	94.6%	25.7%
FP rate	0.0%	0.0%	0.0%	46.7%	3.5%	43.2%
<i>unencrypted</i> MP3						
TP rate	0.0%	0.0%	0.0%	63.8%	95.0%	31.2%
FP rate	0.0%	0.0%	0.0%	32.3%	0.2%	32.1%
<i>encrypted</i> MP3						
TP rate	0.0%	0.0%	0.0%	58.6%	96.1%	37.5%
FP rate	0.0%	0.0%	0.0%	48.3%	0.2%	48.1%
<i>unencrypted</i> ZIP						
TP rate	0.0%	0.0%	0.0%	60.0%	90.4%	30.4%
FP rate	0.0%	0.0%	0.0%	29.9%	8.3%	21.6%
<i>encrypted</i> ZIP						
TP rate	0.0%	0.0%	0.0%	55.5%	90.6%	35.1%
FP rate	0.0%	0.0%	0.0%	28.0%	8.9%	19.1%
<i>unencrypted</i> EXE						
TP rate	2.7%	1.3%	0.1%	54.1%	84.9%	30.8%
FP rate	0.0%	0.0%	0.0%	47.3%	16.7%	10.6%
<i>encrypted</i> EXE						
TP rate	0.0%	0.0%	0.0%	56.1%	84.5%	28.4%
FP rate	0.0%	0.0%	0.0%	54.3%	17.2%	37.1%
<i>unencrypted</i> PDF						
TP rate	5.2%	2.5%	3.6%	75.4 %	84.5%	9.1%
FP rate	0.0%	0.0%	0.0%	46.8%	31.8%	15.0%
<i>encrypted</i> PDF						
TP rate	0.0%	0.0%	0.0%	63.2%	84.8%	21.6%
FP rate	0.0%	0.0%	0.0%	45.5%	31.9%	13.6%
<i>unencrypted</i> DOC						
TP rate	0.1%	0.0%	0.0%	65.6%	66.3%	0.7%
FP rate	0.0%	0.0%	0.0%	48.8%	29.2%	19.6%
<i>encrypted</i> DOC						
TP rate	0.0%	0.0%	0.0%	57.6%	67.7%	10.1%
FP rate	0.0%	0.0%	0.0%	46.2%	31.4%	14.8%

491 triggered without requiring any additional software. Therefore, most contem-
 492 porary malware are bundled in executables files. It is for the same reason that

493 email filters and network firewalls generally remove executable attachments.
494 Similarly, we infer that high detection rates for PDF files are due to the small
495 sizes of these files. We argue that antivirus scanners generally scan a small
496 portion of a file’s payload. Infection is detected only if a known malware sig-
497 nature is present in that portion of the file. It can be said to evade detection
498 by COTS antivirus malicious payload should be embedded towards the end
499 of the file. The false positive rates for these AV products are also 0% because
500 they mostly use signature-based scanning techniques.

501 Mahanalobis n -gram detector performs much better than COTS AV software.
502 However, its performance also degrades for the encrypted dataset. In compari-
503 son, our proposed Markov n -gram detector achieves the best average detection
504 rate with a reasonable average false positive rate. Furthermore, the accuracy
505 of the Markov n -gram detector persists for the encrypted dataset. This is be-
506 cause entropy of a random variable is not dependent on the values or image of
507 the underlying random variable. Therefore, a shift in the histogram does not
508 change the entropy rate values.

509 The reason for less than 100% detection rate of Markov n -gram detector can
510 be traced back to our earlier comment: the lower bound on size of detectable
511 embedded malware is roughly set by the block’s size. Now recall that we have
512 used a block size of 1000 bytes, while the size of 23.6% files in the VX Heavens
513 malware dataset ([14]) is less than 1000 bytes. Nevertheless, even under this
514 block’s size limitation, the smallest malware that the proposed Markov n -gram
515 detector is able to detect is `Worm.Win32.Netsp`, which is only 343 bytes. We
516 also note that 8.2% of files in the VX Heavens malware dataset are smaller
517 than 343 bytes. These file comprise malware that the proposed detector is
518 unable to detect.

519 We argue that it is not entirely fair to compare the accuracy of the Markov n -
520 gram detector with the scheme proposed in [2] because their scheme focuses on
521 *detection* while ignoring *location identification*, whereas the detector proposed
522 in this paper caters for both of these objectives. Despite this fact, the accuracy
523 of the Markov n -gram detector is significantly higher than the Mahanalobis
524 detector. The bold right column in Table 5 gives the percentage improvement
525 in the TP rate and the FP rate for the Markov n -gram detector as compared
526 to the Mahanalobis n -gram detector. It can be observed that the TP rate of
527 the Markov n -gram detector is on the average 20.2% and 26.4% greater than
528 the TP rate of the Mahanalobis detector for non-encrypted and encrypted
529 datasets, respectively. Similarly, the FP rate of the Markov n -gram detector
530 is on the average 21.9% and 29.3% smaller than the FP rate of Mahanalobis
531 detector for non-encrypted and encrypted datasets, respectively. This clearly
532 indicates the superior detection accuracy and robustness of our proposed de-
533 tector as compared to the detector proposed by the authors in [2].

534 We, however, do admit that the FP rates for DOC and PDF files are still sig-
535 nificantly high albeit much smaller as compared to the Mahanalobis detector.
536 As our earlier investigation revealed that *embedded objects* are allowed both
537 in PDF and DOC files. The entropy rate at the location of these objects, at
538 times, also shows a significant perturbation. As a result, our detector is mis-
539 lead to classify these objects as malware. We will shortly introduce our hybrid
540 strategy that will solve this problem of high FP rate.

541 Another important conclusion of the research by the authors in [2] is: if the
542 size of the benign file in which the infection is inserted is between 10 KB
543 and 10 MB then on the average the false positive rate of their scheme surges
544 to 50%. In comparison, our scheme has two desirable features: 1) capability
545 to identify block/blocks of benign file in which the infection was inserted; 2)
546 a significantly smaller false positive rate relative to the Mahanalobis n -gram
547 detector. Here, we must note that sizes of more than 90% of the benign files in
548 our dataset also lie in the 10 KB to 10 MB range (average size = 2 MB). This
549 encouraging performance clearly substantiates the potential of the Markov
550 n -gram detector for embedded malware detection.

551 6 Limitations of the Markov n -gram detector

552 In this section, we present the limitations of the Markov n -gram detector
553 proposed in this paper.

- 554 • The first, and perhaps the biggest, shortcoming of the proposed Markov
555 n -gram detector is its high false positive rate for certain types of files. We,
556 however, believe that these false positives can be significantly reduced if
557 we use the proposed detector as a preprocessor to the signature-based de-
558 tector. During this preprocessing stage, the Markov n -gram detector can be
559 utilized to detect the presence and the location of embedded malware inside
560 a benign file, albeit with false positives. We can then extract a small por-
561 tion of the file around the infected location into a separate standalone file.
562 The signature-based detector can then scan the new (extracted) file for the
563 presence of a known malcode signature. Clearly, this hybrid detection strat-
564 egy can significantly lower the false positives, while still maintaining the
565 high detection rate of our detector. This hybrid detection strategy is only
566 realizable because our detector can identify the location of the embedded
567 malware. Our attempts to use a hybrid of COTS AV software and Markov
568 n -gram detector were not successful. Our investigation shows that the ex-
569 tracted portion of the file still remains an embedded malware for COTS
570 AV software because perfect alignment of the extracted portion cannot be
571 achieved due to naivety of our extraction algorithm. This problem can be
572 solved in two ways: 1) development of intelligent extraction algorithms to

573 embark *exact* location of embedded malware so it can be detected by COTS
574 AV software, and 2) development of an in-house signature based detector
575 enabled with deep inspection. Second proposition is not in the scope of cur-
576 rent work, however, we plan to pursue the first proposition as our future
577 work.

- 578 • One form of the embedded malware discussed in this paper is dormant (see
579 Section 3 for more details), which does not pose any direct threat to the
580 victim machine. The dormant form of embedded malware requires another
581 program (such as a trojan) to extract and activate it. As a results, this
582 problem is similar to detecting watermarks or steganographic content. It
583 is unlikely that our scheme will detect a malware embedded using the ad-
584 vanced steganographic embedding schemes. However, the use of advanced
585 steganographic schemes have two major drawbacks when considered in the
586 context of embedded malware: 1) the steganographic embedding and extrac-
587 tion algorithms have high memory and computational overheads; 2) they
588 are media specific, i.e. they are specific for images, audio or video content,
589 so they cannot be generalized to all filetypes. The first drawback implies
590 that the steganographic extraction algorithm should be present on the vic-
591 tim machine. Transfer of extraction program to the victim machine is an
592 additional and undesirable overhead. Also, the computational overhead, at
593 the victim machine, imposed due to the extraction algorithm allows for
594 host based anomaly detection. The second drawback limits the scope of the
595 threat posed by embedded malware.
- 596 • Since the underlying principle of our proposed detector is based on statistical
597 analysis, a crafty attacker may launch a mimicry attack [20] by modifying
598 the malcode to have a benign looking statistical distribution [3]. Polymor-
599 phic attack engines can be used to modify the statistical distribution of a
600 code segment to avoid detection by our proposed detector. This unfortunate
601 limitation is not specific to the Markov n -gram detector and is applicable
602 to any anomaly detector, including the Mahanalobis n -gram detector and
603 COTS AV software [20].

604 7 Related Work

605 A significant amount of research effort has recently been focused towards mal-
606 ware detection. To maintain focus and logical flow of thought, in this section,
607 we describe only the approaches that target embedded malware detection or
608 n -gram based malware detection.

- 609 • Stolfo et al. extended their previous work on identification of filetypes us-
610 ing n -gram analysis in [2]. In their earlier analysis, called fileprint analysis,
611 they calculated 1-gram byte distribution of a file and compared it to vari-
612 ous models of different filetypes for eventual identification of the filetype. In

613 the context of malware detection, their work focused on embedded malware
614 detection only in PDF and DOC files. They used 3 different models for repre-
615 senting the benign distributions namely single centroid, multi-centroids and
616 exemplar files as centroids. Mahalanobis distance was calculated between
617 the distributions obtained from these models and the n -gram distribution
618 of a given file. To avoid repetition, details of these techniques will be pro-
619 vided in subsequent sections.

620 The authors experimented with 1-gram (byte level) and 2-gram (word
621 level) distributions. They tested their proposed scheme on a dataset com-
622 prising 31 benign application executables, 331 benign executables in the
623 System32 folder and 571 viruses. The results of their experiments demon-
624 strated that their scheme was able to detect a considerable proportion of
625 the malicious files. However their approach was not capable of identifying
626 the exact location of the embedded malware in a benign file. Therefore, it is
627 impossible to devise an effective healing strategy for the infected files using
628 their approach.

- 629 • In [3], the authors proposed two approaches for embedded malware detec-
630 tion in Microsoft Word documents. The first approach is based on static
631 analysis and the second approach is based on run time dynamic analysis.
632 In the static analysis approach, they used an open source application to
633 decompose Word files into their constituent structures. They used a 5-gram
634 model for benign and malicious documents because it provided reasonable
635 memory and detection accuracy. Based on the 5-gram model for benign
636 and malicious word documents, a “similarity” score was generated for both
637 models for eventual classification. In dynamic analysis approach, they have
638 employed sandbox-based tests to check OS crashes, unexpected changes to
639 the underlying environment, and nonfatal application errors. However, it
640 is acknowledged by the authors that the dynamic analysis approach is not
641 practical to be used as an independent detection scheme.
- 642 • Kolter et al. present a hybrid approach to detect malicious executables in the
643 wild [8]. They use n -gram analysis to extract features from 1971 benign and
644 1651 malicious samples. The most relevant n -grams (4-grams in the paper)
645 were selected for prediction by computing the information gain. A variety of
646 inductive learning methods were used, such as naive Bayes, decision trees,
647 support vector machines and boosting. Their results show that the boosted
648 decision trees outperform rest of the classifiers.
- 649 • In [9], the authors have proposed a hybrid approach for detection of new
650 malicious code. They use n -gram analysis to automatically generate virus
651 signatures from a given sample of malware and benign programs. Authors
652 have shown that the n -gram signatures are more robust as compared to the
653 traditional signature based schemes.
- 654 • Henchiri et al. propose a data mining approach to generate robust signatures
655 for computer viruses [10]. They claim to find robust signatures using n -gram
656 analysis. They carried out their analysis on different families of viruses. A
657 subset of most frequent n -grams were chosen from each virus family. Feature

658 elimination was carried out on the collected n-grams to get the final set of
659 features that were the eventual signatures. Authors used various machine
660 learning algorithms on the top of these features for final classification. They
661 carried out their tests on a dataset of 1512 viruses and 1488 benign files and
662 achieved up to 93% overall accuracy.

663 8 Conclusions

664 In this paper, we have proposed a novel embedded malware detection scheme
665 based on the principles of statistical anomaly detection. This scheme, to the
666 best of our knowledge, is the first anomaly based malware detection approach
667 that has the capability to locate the position of the infection in an infected file.
668 Our proposed Markov n -gram detector has significantly better detection rate
669 than exiting detectors. Moreover, due to its ability to identify the location
670 of an embedded malware, the proposed detector can provide very low false
671 positive rates when appropriately used in conjunction with a signature-based
672 detector.

673 Acknowledgments.

674 We acknowledge the help of Ch. Junaid Anwar in dataset generation.

675 This work is supported by the National ICT R&D Fund, Ministry of Informa-
676 tion Technology, Government of Pakistan. The information, data, comments,
677 and views detailed herein may not necessarily reflect the endorsements of views
678 of the National ICT R&D Fund.

679 References

- 680 [1] Symantec Internet Security Threat Report XI: Trends for July – December
681 2007, September 2007.
- 682 [2] S. J. Stolfo, K. Wang and W. J. Li, Towards Stealthy Malware Detection,
683 *Advances in Information Security*, Volume 27, pp. 231-249, Springer US, 2007.
- 684 [3] W. J. Li, S. J. Stolfo, A. Stavrou, E.Androulaki and A. D. Keromytis, *A*
685 *Study of Malcode-Bearing Documents*, Detection of Intrusions and Malware &
686 Vulnerability Assessment (DIMVA), Volume 4579 of Lecture Notes in Computer
687 Science, pp. 231-250, Springer Verlag, Zurich, Switzerland, 2007.

- 688 [4] John Leyden, *Trojan exploits unpatched Word vulnerability* The Register, May
689 2006.
- 690 [5] Joris Evers, *Zero-day attacks continue to hit Microsoft*, News.com, September
691 2006.
- 692 [6] David Kierznowski, *Backdooring PDF Files*, September 2006.
- 693 [7] M. Zubair Shafiq, S. A. Khayam, M. Farooq, *Embedded Malware Detection*
694 *using Markov n-grams*, Detection of Intrusions and Malware & Vulnerability
695 Assessment (DIMVA), Volume 5137 of Lecture Notes in Computer Science, pp.
696 88-107, Springer Verlag, Paris, France, 2008.
- 697 [8] J. Z. Kolter, M. A. Maloof, *Learning to Detect and Classify Malicious*
698 *Executables in the Wild*, Journal of Machine Learning Research, Volume 7, pp.
699 2721-2744, 2006.
- 700 [9] T. Abou-Assaleh, N. Cercone, V. Keselj, and Ray Sweidan, *Detection of*
701 *New Malicious Code Using N-grams Signatures*, Advances in Soft Computing
702 Journal, Springer-Verlag, 2003.
- 703 [10] O. Henchiri and N. Japkowicz, *A Feature Selection and Evaluation Scheme for*
704 *Computer Virus Detection*, IEEE International Conference on Data Mining,
705 2006.
- 706 [11] S. Friedl, *Steve Friedl's Unixwiz.net Tech Tips: Analysis of the new 'Code Red*
707 *II' Variant*, Webpage, <http://www.unixwiz.net/techtips/CodeRedII.html>.
- 708 [12] A. Lakhina, M. Crovella, and C. Diot, *Mining anomalies using traffic feature*
709 *distributions*, ACM SIGCOMM, September 2005.
- 710 [13] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, Wiley-
711 Interscience, June 1991.
- 712 [14] *VX Heavens Virus Collection*, VX Heavens, <http://hvx.netlux.org/v1.php>
- 713 [15] *McAfee Anti-virus and Internet security*, McAfee, CA.
- 714 [16] *AVG Anti-Virus and Internet Security*, GRISOFT Inc., FL, USA.
- 715 [17] *Kaspersky Antivirus*, Kaspersky Lab HQ, Moscow, Russia.
- 716 [18] Yinrong Huang, *Vulnerabilities in Portable Executable (PE) File Format For*
717 *Win32 Architecture*, TR, Exurity Inc., Canada, 2006.
- 718 [19] T. Fawcett, *ROC Graphs: Notes and Practical Considerations for Researchers*,
719 TR, HP Labs, CA, 2003-4, USA.
- 720 [20] D. Wagner and P. Soto, *Mimicry Attacks on Host-Based Intrusion Detection*
721 *Systems*, ACM CCS, Nov. 2002.