

Behavioral detection of malware: from a survey towards an established taxonomy

Grégoire Jacob · Hervé Debar · Eric Filiol

Received: 15 June 2007 / Revised: 27 October 2007 / Accepted: 31 January 2008 / Published online: 21 February 2008
© Springer-Verlag France 2008

Abstract Behavioral detection differs from appearance detection in that it identifies the actions performed by the malware rather than syntactic markers. Identifying these malicious actions and interpreting their final purpose is a complex reasoning process. This paper draws up a survey of the different reasoning techniques deployed among the behavioral detectors. These detectors have been classified according to a new taxonomy introduced inside the paper. Strongly inspired from the domain of program testing, this taxonomy divides the behavioral detectors into two main families: simulation-based and formal detectors. Inside these families, ramifications are then derived according to the data collection mechanisms, the data interpretation, the adopted model and its generation, and the decision support.

1 Introduction

Even though behavioral detection seems a recent trend, in antivirus products as well as in virology research, its principles are not really new. In 1986, Cohen [1,2] already established a basis for behavioral detection within his first formal works. He made his point that viruses, just like any other running program, use the services provided by the system. Predicting the viral nature of a program by its behavior was

then equivalent to defining what is, and what is not a legitimate use of the system services. This problem was eventually reduced to an appearance analysis of the inputs sent to the system, which is undecidable. Basically, this definition is strongly linked to the operating system but it can easily be extended to the use of any hardware or software resource (processor, memory, programs). This extended definition is often referred as function-based detection. The difference remains a question of perimeter explaining that function-based and behavioral detections are considered indifferently along the article.

1.1 Two opposite approaches for behavioral detection

As stated by Cohen, two opposite approaches can apprehend the problem of behavioral detection. The first approach is to model the behavior of legitimate programs and measure deviations from this reference. The great advantage of this approach lies in its capacity to detect completely unknown viral strains. Nevertheless, defining a global behavior for programs reveals itself extraordinary complex. An obvious reason is the multitude of applications with different natures existing on a system. A web or mail client exhibits an intensive use of the network facilities whereas a multimedia player decodes large buffers of data and renders them over physical devices such as the graphic or sound cards. No common characteristics can be extracted and a different profile is required for each kind of application. Moreover the available information is too important for each program (several megabytes of code, hundreds of system calls) to be considered wholly. As a consequence, legitimate models are always statistical, thus prone to false positive and non resilient to major environment changes. It explains why, in virology, the second opposite approach of modelling and detecting suspicious behaviors is mainly adopted. When a set proves too

G. Jacob (✉) · H. Debar
France Télécom R&D, Caen, France
e-mail: gregoire.jacob@orange-ftgroup.com;
gregoire.jacob@gmail.com

H. Debar
e-mail: herve.debar@orange-ftgroup.com

G. Jacob · E. Filiol
French Army Signals Academy,
Virology and Cryptology Lab, Rennes, France
e-mail: eric.filiol@esat.terre.defense.gouv.fr

complex to be defined exhaustively, the problem can intuitively be addressed working on its complementary. The main drawback is that unknown malware can no longer be detected as soon as they use innovative viral techniques.

It is interesting to parallel antivirus products with intrusion detection systems, where the perception is diametrically opposed. In the intrusion domain, behavioral detection is based on legitimate models whereas the suspicious models used in virology are considered as simple signatures for knowledge-based detection, also called misuse detection [3,4]. Modelling legitimate behaviors goes back to the early works on intrusion detection published by Anderson [5] and Denning [6]. It still remains an active research field as it is clearly impossible to generate misuse signatures for the thousands of vulnerabilities discovered every year. Such models are out of the scope of this paper but, for further information, the reader is invited to refer to the works of Forrest et al. [7] on host-based intrusion detection and the works of Zanero [8] on the use of Markovian Models to capture legitimate uses of systems. To go back to our main focus, viral techniques are less numerous than vulnerabilities and misuse models seem more adequate to the present problem of malware detection.

1.2 Paper contribution and organization

In virology, the domain of behavioral detection shows an increasing activity both in commercial products and research. Paradoxically, no global survey covering this domain has been published. A striking multitude of behavioral detection systems can be observed without any will of consistency in the used vocabulary and designations. To our knowledge, this taxonomy dedicated to behavioral detection in virology is the first of its kind, contrary to intrusion detection where the literature is abounding. The novelty of our approach lies in the parallel made with the domain of program testing which makes a distinction between simulation-based and formal verifications. The domain of behavioral detection should benefit greatly from a consistent reference taxonomy. In effect, this taxonomy should remove the divisions between the different sub-domains of behavioral detection, helping information sharing and reuse.

The scope of this taxonomy has been defined as wide as possible, according to the definition of behavioral detection given in introduction. The virology point of view has been willingly chosen, meaning that the modelling of suspicious behaviors has been implicitly considered. As the need arises, relevant intrusion detection papers are also given as additional references. Globally, the paper has been organized as follows: Sect. 2 explains the recent interest in behavioral detection by the predicted failure of appearance detection, Sect. 3 describes a generic behavioral detection system, Sect. 4 is the core part of the article introducing the taxonomy,

and Sect. 5 gives an illustrative overview of both existing commercial products and research prototypes.

2 Why behavioral detection may succeed where form-based detection will undeniably fail

Historically, appearance detection also called form-based detection has been the first technique used to fight malware and still remains at the heart of nowadays antivirus software. These detection techniques search system objects such as files for suspicious byte patterns referenced in a base of signatures. These betraying patterns must exhibit a discriminating character combined with non-incriminating properties for legitimate programs [9, p. 147]. Even if these purely syntactic signatures can precisely identify the threat and name it, form-based techniques are bound to detect known malware or trivial variants.

On the opposite, behavioral signatures are no longer simple byte patterns but complex meta-structures carrying dynamic aspects and a semantic interpretation. Programs with distinct syntaxes can basically have an identical behavior captured by a single behavioral signature. As a consequence, a behavioral signature no longer identifies a single piece of malware but a whole class of malware. Behavioral detection is thus more generic and more resilient to modifications than form-based detection. On the other hand, a precise identification of a piece of malware inside its class is no longer possible, which can be problematic when choosing the relevant countermeasure. Nevertheless, behavioral detection should bring a solution to two of the major problems encountered by form-based detection.

2.1 The signature extraction problem

Form-based detection provides undeniable advantages for operational use. It uses optimized pattern matching algorithms with controlled complexity and very low false positive rates. Unfortunately, form-based detection proves completely overwhelmed by the quick evolution of the viral attacks. The bottleneck in the detection process lies in the signature generation and the distribution process following the discovery of new malware.

The signature generation is often a manual process requiring a tight code analysis that is extremely time consuming. Once generated, it must be distributed to the potential targets. In the best cases, this distribution is automatic but if this update is manually triggered by the user, it can still take days. In a context where worms such as Sapphire are able to infect more than 90% of the vulnerable machines in less than 10 min, attacks and protection do not act on the same time scale.