Provably Secure Virus Detection: Using The Observer Effect Against Malware*

Richard J. Lipton¹, Rafail Ostrovsky^{†‡2}, and Vassilis Zikas^{§‡3}

- 1 Department of Computer Science, Georgia Institute of Technology, Atlanta, USA
 - rjl@cc.gatech.edu
- 2 Department of Computer Science & Department of Mathematics, University of California, Los Angeles, USA rafail@cs.ucla.edu
- 3 Department of Computer Science, Rensselaer Polytechnic Institute, Troy, USA vzikas@cs.rpi.edu

Abstract

Protecting software from malware injection is one of the biggest challenges of modern computer science. Despite intensive efforts by the scientific and engineering community, the number of successful attacks continues to increase.

This work sets first footsteps towards a provably secure investigation of malware detection. We provide a formal model and cryptographic security definitions of attestation for systems with dynamic memory, and suggest novel provably secure attestation schemes. The key idea underlying our schemes is to use the very insertion of the malware itself to allow for the systems to detect it. This is, in our opinion, close in spirit to the quantum Observer Effect. The attackers, no matter how clever, no matter when they insert their malware, change the state of the system they are attacking. This fundamental idea can be a game changer. And our system does not rely on heuristics; instead, our scheme enjoys the unique property that it is proved secure in a formal and precise mathematical sense and with minimal and realistic CPU modification achieves strong provable security guarantees. We envision such systems with a formal mathematical security treatment as a venue for new directions in software protection.

1998 ACM Subject Classification D.4.6 Security and Protection (K.6.5)

Keywords and phrases Cryptography, Software Attestation, Provable Security

Digital Object Identifier 10.4230/LIPIcs.ICALP.2016.32

1 Introduction

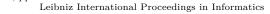
Protecting software from malware injection is a major goal of computer security. Nonetheless, the problem of basing solutions on strong theoretical foundations does not seem to have received sufficient attention in the malware protection literature. In this work we suggest a novel provably secure and practically efficient paradigm for software protection against

[§] This author was supported in part by the Swiss NSF Ambizione grant PZ00P2_142549.



43rd International Colloquium on Automata, Languages, and Programming (ICALP 2016). Editors: Ioannis Chatzigiannakis, Michael Mitzenmacher, Yuval Rabani, and Davide Sangiorgi; Article No. 32; pp. 32:1–32:14





^{*} A preliminary version of this work was presented at the Securing Computation Workshop, Simons Institute, UC Berkeley. US Patent Pending [22].

 $^{^\}dagger$ This author was supported in part by the NSF award CNS-1136174.

[‡] Work done in part while the author was visiting the Simons Institute for the Theory of Computing, supported by the Simons Foundation and by the DIMACS/Simons Collaboration in Cryptography through NSF grant #CNS-1523467.

32:2 Provably Secure Virus Detection

arbitrary malicious code injection. We provide the first formal cryptographic model tailored to detect malware injection in modern computers, which we envision as the basis for fruitful research on provable secure detection and prevention of malware leading into the next generation of secure systems. Importantly, we provide first matching solutions that are accompanied by mathematical proofs that any memory oblivious injection will be caught with arbitrary high probability. Our solutions are practically efficient, demonstrating that provable security does not come at a too high price.

Our proposed system has the following desirable properties: 1) it requires minor to no changes to the CPU specification to provably defend against injection attacks that can not read the code before they inject their malware; 2) it only affects the performance of the program by a modest amount. The first property implies that it can be used on today's computers. The second property means that it can be practical, since performance is critical in most applications; this is obvious even in its simplified form presented here, where various optimizations have been avoided for the sake of clarity in the presentation. Finally, we allow attacks both spatial and temporal freedom: the attack can occur at any time during the execution and on any part of the memory. It can attack code as well as data.

1.1 Overview of our Contributions

We put forward and formally define the notion of a virus detection scheme (in short, VDS) which compiles any given program W (and its data) into a new secured program \widetilde{W} that performs the same computation as W but allows us to detect any virus injected in the memory at any point of the execution path. The detection is done via a provably secure challenge-response mechanism between the machine executing the (compiled) software and a verifying external device. Importantly, we insist that the verification algorithm be very simple, and in particular that it be executed by a very lightweight device (Our constructions require that the verifier only does an encryption and compares strings for equality.) Thus, the role of the verifier can be, for example, played by the user with a smartphone, or by a compact and simple, intrusion-free hardware module that could even be part of the CPU.\(^1\) Moreover, the verification uses public key techniques, where the verifier needs only the public key. This restricts the attack surface to the machine executing the secured code, as compromising the verifier's privacy gives an attacker no advantage in breaking our scheme.

We provide our formal cryptographic definitions of VDS security based on the well-studied Random Access Machine (RAM) model, slightly adapted to be closer to an abstract version of a modern computer following the *von Neumann architecture*. The suggested model corresponds to a simplistic *closed-system* abstraction of software execution: The code to be executed along with its associated data is loaded initially on the random access memory (RMEM,) which, through the execution, only communicates with the CPU.²

In more detail, a secure VDS is described as follows. The software to be executed is compiled, prior to being loaded onto the memory, to a secure version. Importantly, this compilation does not necessarily need the source code but for a wide class of software (non-self-modifying code) can be applied on the binary (machine-language) code of the program; thus, the compilation can be performed by the program vendor (this might allow for further

¹ Unlike software-based attestation techniques (cf. Section 1.2), the verification in our scheme can be done even over an insecure network, e.g. the Internet.

² In particular, the current theoretical model does not address how the data is exchanged with secondary storage devices, e.g., flash memory. It is part of ongoing research to extend this model to more complicated architectures.

efficiency optimizations) or by the user without requiring any reverse engineering. The compiler, in addition to the (binary code) of the program and its data, uses a randomly chosen *compilation key* K.

To check/verify that the software running on a system has not been attacked by a virus (or to detect such an attack in case it has occurred), the following detection process is executed. The user issues a challenge c that depends on the compilation key K and the system is required to produce a reply which passes a specified verification test. As already discussed, we arrange things in such a way that the verifier does not need to know the compilation key, and he can just know some public information on it. Formally, this is done by the use of public key cryptography, where the compilation key is the secret key, and the information held by the verifier, which we refer to as the verification key, is the corresponding public key. The security of the VDS requires that if the RAM has not been attacked then it can always reply to the challenge in an accepting manner (verification correctness); otherwise, any reply will be rejected with high probability (security).

We provide concrete instantiations of practical VDSs which, depending on the assumptions about the CPU they are executed on, achieve from a reasonably strong security (namely protection against continuous injection of moderate-sized virus) to security against arbitrary small viruses injected on locations that might depend on key-share locations. Informally, we prove the following result for our VDSs.

▶ Theorem (informal). Under standard complexity assumptions our VDSs compile any non-self-modifying software into a secure version that detects any malware injection with very high probability.

Our schemes are independent of the virus code, are platform-independent, and might require only small modifications to the common CPU architecture (in particular, they do not assume tamper-resilient hardware); and, with appropriate optimizations we can make it that they only affect the performance of the executed software by a practically acceptable amount. As a useful side effect, our scheme is even able to detect hardware errors, e.g., random bit-flips in the memory; as demonstrated in [5], such errors might have important consequences on the security of the executed software.

To construct our VDS we start by constructing a scheme satisfying a weak(er) notion of security, which (1) catches only injections of continuous and sufficiently long viruses, and (2) assumes that the response is computed by an external (trusted) device that is given access to the state of the attacked system. We give a formal definition of this weaker notion which depending on the application and the desired security level might already be satisfactory. We then proceed and gradually modify (strengthen) our scheme. The first modification ensures that the compiled program can compute the response by itself, i.e., without the help of an external trusted device. We note that the security of our weakest scheme is information theoretic (it is based on the perfect privacy of one-time pad encryption), whereas the proof of the more secure scheme requires a leakage-resilient encryption algorithm [9].

The second modification uses a simple message authentication code (in short, MAC) to remove the limitation to sufficiently long and continuous virus injection. A MAC authenticates a value using a random key, so that an attacker without access to the key cannot change the value without being detected. More concretely, to authenticate a word w our scheme relies on the standard MAC [16] that uses two keys K_1 and K_2 , and compute a MAC tag as $t = w \cdot K_1 + K_2$, where w, K_1 , and K_2 are interpreted as elements of an appropriate large arithmetic field. Observe that such arithmetic operations are implemented in hardware in existing CPUs. We assume, however, that the executing CPU has a slightly extended instruction-set which, informally, has each "load" instruction, i.e., read-from-memory instruction, check that the

32:4 Provably Secure Virus Detection

MAC is correct before it loads the word on the CPU registers. We stress that this only needs the architecture to provide us with a very simple extra piece of microcode. By engineering the CPU's hardware in a smart way so that these operations are done on the circuit level, we expect to be able to reduce the effect of our compiler to a barely noticeable slowdown. Fortunately, the new architecture that Intel recently announced promises to embed such microcode in its next-generation microchips [24].

Overview of our VDS Construction. So how can we detect an intrusion that we have never seen before? How can we arrange that any injection of malware into our program by an attacker who does not know the compilation-key will be detected? The answer is that we will hide the compilation-key/secret in our program in a way that ensures that with very high probability, any injection by an adversary must destroy the secret. Once this secret is destroyed, the adversary may indeed be able to take over the program, but will be quickly detected. We note that current systems may fail to detect such an attack forever – for the type of injection we protect against, we are able to detect it in seconds.

A point: We can imagine situations where even such a swift detection of an attack is not sufficient to stop potential harm. An attacker could, in some situations, do immense damage even if he is in control of a program for only a fraction of a second. Furthermore, our approach is not targeted towards fixing buggy software or detecting malware-including software which might be (unknowingly) installed by the user. Rather, our method offers protection against "unauthorized" injection of malicious code, e.g., direct injection to the memory or radiation attacks. But we view the ability to provably detect any such attack, new, old, clever, or not, as a big improvement over the current state of the art.

So how can we do this? Let W be a program that we wish to protect from malware. We recompile W to \widetilde{W} . The idea is that W and \widetilde{W} must compute the same thing, run in about the same time, and yet \widetilde{W} must be able to detect itself against the insertion of malware. An obvious idea is to have \widetilde{W} periodically check to see if its code or data have been maliciously changed. The trouble with this approach is that the attacker could easily inject his own code, take over the checker, and thereby disable the checking. In short: who checks the checker?

Here is how we proceed. The protected program W operates normally most of the time. Periodically it is challenged by another machine to prove that it has a secret key K. This key is known only to \widetilde{W} and not to the attacker. If \widetilde{W} has not been attacked, then it simply uses the key K to answer the challenge and thus proves that it is still operating properly.

The issue is what happens if \widetilde{W} has been attacked and some code has been injected into it. We can arrange \widetilde{W} so that no matter how the injection of code has happened the key K is lost. The attacker's very injection will have changed the state of \widetilde{W} so that it is now in a state that no longer knows the key K. This is the analog of the quantum observer effect, often referred to as the Heisenberg Uncertainty Principle: the attacker has damaged the state of \widetilde{W} so that the key has been destroyed.

Making the attack destroy the key is the central idea here. If \widetilde{W} simply stores K somewhere in memory, the attacker will take over the program, look around for the key, and likely find the key;³ this defeats the protection, since now it can answer the periodic challenges just like \widetilde{W} could. To resolve the above issue we distribute the key K all through memory by using what is called secret sharing [30]. This is a standard method in cryptography that breaks a small key, like K, into many small pieces, called shares. These pieces are cleverly placed throughout all of \widetilde{W} 's memory. Our secret sharing allows K to be reconstructed only

³ Note that we do not assume private memory, i.e., memory which the CPU cannot read.

if all the pieces are left untouched – if any are changed in any way, then it is impossible to reconstruct the key. Obviously, if there has been no attack, then in normal operation \widetilde{W} can reconstruct the key from the pieces and answer the challenges. However, if the pieces are not all intact, then \widetilde{W} will be unable to answer the next challenge.

This is the high-level idea of our method: hide the key via secret sharing, and rely on the attacker to destroy at least one share of the key. This destruction is irreversible and makes the attack fail the next challenge.

There is one additional point that we must mention. We must arrange that \widetilde{W} can be attacked at anytime, including when the system has collected all the shares and reconstructed the key K on its CPU. This is a very dangerous time, since if \widetilde{W} is attacked at this moment the fact that shares are destroyed does not matter. The attacker can simply use the reconstructed key K. We avoid such time-of-check-time-of-use (TOCTOU) attacks [25] by actually using two keys K_1 and K_2 in tandem. Both are stored via secret sharing as before, but now one must have both in order to answer the challenges. We arrange that \widetilde{W} only reconstructs one key at a time and this means that an attacker must destroy either or both of the keys. This handles the above dangerous situation and makes our method provably secure.

Tying up the last loose end, we treat the case of very small viruses – i.e., ones that consist of a single bit or just a few bits – and viruses that know (or guess) the exact memory locations of key-shares in advance and therefore might not need to overwrite them. We armor our system to defend even against such tiny/informed viruses by employing an additional lightweight cryptographic primitive, namely a MAC. We use the MACs in a straightforward and efficient manner: we authenticate each word in W using the key-shares (also) as MAC keys, and require that for any word that is loaded from the random access memory, the CPU verifies the corresponding MAC before further processing the word. Thus, if the MAC does not check out we detect it (and immediately destroy the share), and if the adversary changes the MAC key, he loses one of the shares of one of the two secrets, and we detect that as well.

We enforce the above checks by assuming a modified CPU architecture (with only very small amount of additional computation, so the CPU power consumption is not affected). More specifically, to get the most out of the MAC functionality we need that when the CPU executes the program, it verifies authenticity of the words it loads from the memory. That is, the load operation of the CPU loads a block of several (in our case, five) words – this is already the case in most modern CPUs – including the program word, the MAC-tag and the corresponding keys, and check with every load instruction that the MAC verifies before further processing the word. Importantly, our MAC uses just one field addition and one field multiplication per authentication/verification; the circuits required to perform these operation are actually trivial compared to the modern CPU complexity and are part of many existing CPU specifications.

But to obtain security against any injection, we need one more trick: instead of using the actual key-shares in the generation/verification of the MAC tag, we use their hash-values. This ensures that the virus cannot manipulate the keys and forge a consistent MAC unless he overwrites a large part of them. We anticipate that given our goal's potential impact on security, such functionality might be included in the next generations of CPUs [24].

1.2 Related Literature

The literature on defense mechanisms against malware-injection attacks is vast. Due to its urgency and importance, the problem has been intensively researched. In what follows we provide a brief overview of the directions most related to our approach.

Closest related to our goals is the literature on verifying the authenticity/integrity of the internal state of computing devices to confirm that they have not been attacked by malware,

32:6 Provably Secure Virus Detection

a mechanism usually referred to as attestation. A rough taxonomy divides this literature in three general categories, namely hardware-based attestation [2, 6, 11, 21], software-based attestation [20, 29, 28, 1, 17, 18], and remote attestation [6, 11, 21, 14, 4]. Due to the similarity in the goals of attestation with our system, we present a detailed overview of the corresponding literature in the full version of this work [23]. As a general note, we stress that similarly to other practical defense mechanisms, attestation schemes do not come with a formal proof of security or even a theoretical security model. In contrast, our scheme is not only backed by a formal mathematical proof, but has several additional advantages with respect to existing attestation schemes, which are highlighted in the full version of this work [23]. We note, however, that much of the attestation literature allows the system to leak its entire state; although our schemes do not account for such leakage (in particular, the key shares should not leak), one can use sharing-refreshing techniques, e.g., [26, 3], to add protection against periodic leakage. The details of such a scheme are subject of future work.

On the most theoretical side of cybersecurity, cryptography provides solutions whose security is backed by rigorous mathematical proofs that typically reduce hardness of breaking the scheme to hardness of solving a mathematical problem, e.g., factoring. But with only a few exceptions, e.g., [7, 13, 10, 19, 15, 27], the cryptographic literature has not, to the best of our knowledge, targeted the problem of malicious code injection. And in contrast to the security-engineering research, cryptographic solutions that do target this problem are often inefficient and/or adopt a too-abstract model of computation which makes them inapplicable or impractical for today's systems.

1.3 Organization of the Paper

In Section 2 we describe the model of computation and the virus injection model, and in Section 3 we provide our security definitions for VDSs. Our VDS constructions are then described in Sections 4 and 5. Due to limited space, some of the formal definitions and proofs have been moved to [23].

2 The Model

In this section we provide an abstract specification of our model of computation. We use as our basis the well known Random Access Machine (RAM) model but slightly adapt it to be closer to an abstract version of a modern computer following the von Neumann architecture. In a nutshell, this modification consist of assuming that both the program and its data are written on the RAM's random access memory⁵ which is polynomially bounded. A RAM \mathcal{R} consist of two components: A Random Access Memory (in short RMEM) and a Central Processing Unit (in short CPU). The RMEM and the CPU communicate in fetch-and-execute cycles, aka CPU cycles, where the number of CPU cycles is the default complexity measure of a RAM. We will refer to a CPU cycle as a round in the RAM execution.

The memory RMEM is modeled as a vector MEM of $m=|\mathtt{MEM}|=poly(k)$ words, where k is the, often implicit, security parameter. Each word is an L-bit string, where we assume that L is linear in k. (Wlog, in the following we will assume that L=k.) For $i\in\{0,1,\ldots,m-1\}$ we denote by $\mathtt{MEM}[i]$ the ith word, i.e., the contents of the i-th RMEM register. The CPU

⁴ An exception here is [4] which uses an idea similar in spirit to ours – i.e., sharing a secret through the memory – for provably protecting against a class of heap overflow attacks (cf. [23] for a comparison).

⁵ In the literature, a RAM with this modification is usually called a *Random Access Stored-Program* machine [12] (in short, RASP).

consists of a much smaller number of L-bit registers – in this work we assume that the total amount of storage of the CPU is linear in the security parameter k – and an instruction set \mathcal{I} that defines which operations can be performed on the registers, and how data are loaded-to/output-from the CPU. The CPU registers include a read-only input-register and an output register which correspond to its interface with its environment (the user). The CPU registers are modeled as an array REG of words, where we denote by REG[i] the content of the ith CPU register. We denote a CPU by the pair $\mathcal{C} = (\text{REG}, \mathcal{I})$ of the vector REG of registers and the instruction set \mathcal{I} . We denote a RAM with CPU \mathcal{C}_k and RMEM MEM $_k$ as $\mathcal{R}_k = (\mathcal{C}_k, \text{MEM}_k)$. The state of \mathcal{R}_k at any point in the protocol execution is the vector (REG $_k$, MEM $_k$) including the current contents of all its CPU and RMEM registers. Details of the model and formal definitions can be found in [23].

To allow for asymptotic security definitions – where the word size, the size of the CPU (i.e., number of its registers), and the size of the memory depend on the security parameter – we often consider a family of RAMs, $\mathcal{R} = \{\mathcal{R}_k\}_{k \in \mathbb{N}}$ with $\mathcal{R}_k = (\mathcal{C}_k = (\mathtt{REG}_k, \mathcal{I}_k), \{\mathtt{MEM}_k\})$. (The size of each word processed by \mathcal{R}_k is L = k.) The RAM families considered in this work have the following property: The instruction set \mathcal{I}_k is the same for all values of the security parameter k. In particular, we assume that all elements of a RAM family, have the same constant number $c = \mathcal{O}(1)$ of CPU registers, where each register of \mathcal{R}_k is of size k, and there is some set of instructions \mathcal{I} that includes operations defined over strings of arbitrary size such that $\mathcal{I}_k = \mathcal{I}$ for every $k \in \mathbb{N}$. (Of course for each value of k, \mathcal{I}_k corresponds to the restriction of \mathcal{I} on k-bit strings.)

Software Execution and Virus Injection. A program to be executed on a RAM $\mathcal{R} = (\mathcal{C}, \text{MEM})$ is described as a vector $W = (w_0, \dots, w_{n-1}) \in (\{0, 1\}^L)^n$ of words that might be instructions, addresses, or program data. To avoid confusion, we refer to such a vector including the (binary of a) software and its corresponding data as $code\ for\ \mathcal{R}$. By convention, whenever, for a RAM family \mathcal{R} we say that W is code for \mathcal{R} , we mean that W is code for the element $\mathcal{R}_k \in \mathcal{R}$ with register size as long as the word size of W and instruction set that includes all the instructions used by W. The execution of a program proceeds as follows: The code W is loaded onto the memory MEM. Unless stated otherwise, we assume that W is loaded sequentially on the first n = |W| locations of MEM, where all remainder locations are filled with (no_op) instructions. The user might give input(s) to \mathcal{R} by writing them on its input register. The RAM starts its execution by fetching the word of the RMEM which the program counter points to, i.e., MEM[pc] (wlog, initially pc = 0). We assume that the RAM is reactive, i.e., any new input written on its input register, makes the RAM resume its computation even if it had halted with output (cf. [23]).

A CPU is *complete* (also referred to as *universal*) if given sufficient (but polynomial) random access memory it can perform any efficient deterministic computation (for a formal definition see [23]). We at times refer to a RAM (family) with a complete CPU as a complete RAM (family).

Modeling Virus Attacks. In our model, a virus attacks a RAM by injecting its code on selected locations of the memory RMEM. More formally, an ℓ -bit virus is modeled as a tuple $\mathbf{v} = (\vec{\alpha}, V) = ((\alpha_0, \dots, \alpha_{\ell-1}), (b_0, \dots, b_{\ell-1}))$, where each $\alpha_i \in \vec{\alpha}$ is a location in the memory and each b_i is a bit. The effect of injecting a virus \mathbf{v} into a RAM $\mathcal{R} = (\mathcal{C}, \text{MEM})$, is to have, for each $\alpha_i \in \vec{\alpha}$, the α_i -th bit on the memory – i.e., the $(\alpha_i \mod L)$ -th bit of the word written

 $^{^{6}}$ In slight abuse of notation we at times drop the security parameter whenever it is clear from the context.

on register $\mathtt{MEM}[\lfloor \frac{\alpha_i}{L} \rfloor]$ – (over)written with b_i . We say that \mathbf{v} is valid for \mathcal{R} if the following properties hold: (1) $\alpha_i \neq \alpha_j$ for every $\alpha_i, \alpha_j \in \vec{\alpha}$, and (2) $\alpha_i \in \{0, \dots, L(|\mathtt{MEM}|-1)\}$ for every $\alpha_i \in \vec{\alpha}$. Furthermore, we say that \mathbf{v} is non-empty if $|\vec{\alpha}| > 0$. Note that in this work we do not consider viruses which inject themselves on the CPU registers.

3 Virus Detection Schemes

We next formalize our notion of provably secure virus detection. More concretely, we introduce the notion of a *virus detection scheme* (in short VDS) which demonstrates how to compile a given program (and its data) into a new program which allows us to detect any virus injection. The detection is done via a challenge-response mechanism.

- ▶ **Definition 1.** A virus detection scheme (VDS) \mathcal{V} consists of five (potentially) randomized algorithms, i.e., $\mathcal{V} = (\text{Gen}, \text{Comp}, \text{Chal}, \text{Resp}, \text{Ver})$, defined as follows:⁷
- Gen is a key-generation algorithm; it computes a pair (compilation-key, verification-key), i.e., $(K_c, K_v) \stackrel{\$}{\leftarrow} \text{Gen.}^8$
- Comp on input the description $\mathcal{R} = \{(\mathcal{C}_k, \mathtt{MEM}_k)\}_{k \in \mathbb{N}}$ of a RAM family, some code W for \mathcal{R}_k , and a compilation key K_c , Comp outputs a new code \widetilde{W} for \mathcal{R}_k (which we will refer to as secure code); i.e., $\widetilde{W} \stackrel{\$}{\leftarrow} \mathtt{Comp}(\mathcal{R}_k, W, K_c)$.
- Chal on input a verification key K_v , and a string $z \in \text{INP}_{\texttt{Chal}} \subseteq \{0,1\}^{poly(k)}$, Chal outputs poly(k)-bit string c called the challenge; i.e., $c \xleftarrow{\$} \texttt{Chal}(z, K_v)$.
- Resp on input a string $c \in \text{OUT}_{\texttt{Chal}}$ and some code \widetilde{W} , Resp outputs a poly(k)-bit string y called the response; i.e., $y \xleftarrow{\$} \text{Resp}(c, \widetilde{W})$.
- Ver on input a verification key K_v , a message $z \in \text{InP}_{Chal}$, a challenge c and a response y, Ver outputs a bit b; i.e., $b \stackrel{\$}{\leftarrow} \text{Ver}(K_v, z, c, y)$. We say that Ver accepts iff b = 1.

A VDS should satisfy the following four security properties (see [23] for formal definitions). The first property is verification correctness which, intuitively, guarantees that if the RAM has not been attacked, then the reply to the challenge is accepting. The second property is compilation correctness, which intuitively ensures that the compiled code \widetilde{W} performs the same computation (on \mathcal{R}) as the original code W. In an application of a VDS, the verifier inputs the challenge at a point of his choice and checks that the reply verifies according to the predicate Ver. Thus, the third property of a VDS is self-responsiveness: the secured code \widetilde{W} includes code that on some special input emulates algorithm Resp on the RAM it executes. Finally, the fourth property requires detection accuracy which states that if some non-empty malware gets injected onto the RAM, then it is detected with high probability. This is one of the most challenging properties to ensure and is the heart of any VDS. We next specify this property by means of a security game between an adversary Adv who aims to inject a virus on a RAM, and a challenger Ch who aims to detect it.

 $^{^{7}}$ All five algorithms below take as an additional input the security parameter k, which is omitted for compactness.

⁸ Note that if we instantiate the VDS with symmetric-key cryptography then $K_c = K_v$.

⁹ This description of the family includes the word-size, the size and addresses of the CPU and RMEM registers and (an encoding) of the instruction set \(\mathcal{I} \).

¹⁰ Without loss of generality, whenever we refer to the secure (i.e., compiled by Comp) version \widetilde{W} of some code W for a RAM family \mathcal{R} we will implicitly assume that \mathcal{R} has enough memory for writing \widetilde{W} on it, i.e., if the word size of \widetilde{W} is k then $|\mathtt{MEM}_k| > |\widetilde{W}|$.

3.1 Detection Accuracy as a Security Game

At a high-level, the security game, denoted by $\mathcal{G}_{\text{VDS}}^{\mathcal{R},\mathcal{V},W}$, proceeds as follows: The challenger Ch runs the key-generation algorithm to obtain a key-pair (K_c,K_v) , and compiles some code W for RAM $\mathcal{R}=(\mathcal{C},\text{MEM})$ onto a new code W for \mathcal{R} by invocation of algorithm Comp; Ch then emulates an execution of W on \mathcal{R} , i.e., emulates its CPU cycles and stores its entire state at any given point. The adversary is allowed to inject a virus of his choice on any location in the memory MEM. Eventually, the challenger executes the (virus) detection procedure: It computes a challenge c by invocation of algorithm $\mathrm{Chal}(z,K_v)$, and then feeds input (check, c) to the emulated RAM and lets it compute the response y. The adversary wins if the output b of the verification algorithm with this response y equals 1. To capture worst case attack scenarios, we allow the adversary to inject his virus at any point during the RAM emulation and make no assumption as to how many rounds the RAM executes after the virus has been injected and before it receives the challenge. Furthermore, we make no assumptions on how much information the adversary holds on the original code W or on the inputs/outputs of \mathcal{R} . The formal description of the security game $\mathcal{G}_{\mathrm{DS}}^{\mathcal{R},\mathcal{V},W}$ can be found in [23].

- ▶ **Definition 2.** We say that a virus detection scheme $\mathcal{V} = (\texttt{Gen}, \texttt{Comp}, \texttt{Chal}, \texttt{Resp}, \texttt{Ver})$ is secure for RAM family \mathcal{R} if it satisfies the following properties:
- 1. V is verification correct, compilation correct, and self-responsive.
- 2. For sufficiently large k for any code W for \mathcal{R}_k and any polynomial adversary Adv in the game $\mathcal{G}_{\text{VDS}}^{\mathcal{R}_k,\mathcal{V},W}$ who injects a valid non-empty virus the following holds: $\Pr[b=1]$ is negligible, where the probability is taken over the random coins of Adv and Ch.

In our constructions we restrict our statements to a certain class of code that satisfies some desirable properties making the compilation easier. We will then say that the corresponding VDS is secure with respect to the given class of code.

The Repeated Detection Game $\tau \mathcal{G}_{\text{VDS}}^{\mathcal{R},\mathcal{V}_1,W}$. Definition 2 requires that the adversary is caught even when he injects its virus while the RAM is executing the Resp algorithm. We next describe a relaxed security game, which provides a useful guarantee for practical purposes. In this game the virus detection (challenge/response) procedure is executed multiple times periodically (on the same compiled code); the requirement is that if the adversary injects its virus to the RAM at round ρ , then he will be caught by the first invocation of the virus detection procedure which starts after round ρ . Note that all executions use the same compiled RAM program and therefore the same key K. The corresponding security game $\tau \mathcal{G}_{\text{VDS}}^{\mathcal{R},\mathcal{V},W}$ which involves τ executions of the detections procedure is detailed in [23]. The security definition is similar to Definition 2 but requires that any virus that is injected in the RAM will be caught in the first virus detection attempt performed after the injection. Our VDSs will be proven secure in this repeated detection game. In [23] we describe a simple trick which transforms our VDSs to ones that are secure in the standard game using a secure hash function.

4 A Software-based VDS for Continuous Moderate-size Virus

In this section we provide a VDS which is secure assuming the virus is linear in the security parameter and that it is injected in continuous memory locations, i.e., if $\mathbf{v} = ((\alpha_0, \dots, \alpha_{|V|-1}), V)$, then $\alpha_i = \alpha_{i-1} + 1$ for all $i \in [|V|-1]$. This captures the entire class of tampering accounted for in [4]. In Section 5 we will show how to get rid of these restrictions. Our construction proceeds in two steps. First, we show how to construct a VDS \mathcal{V}_1 which

achieves a weaker notion of security that, roughly, does not have self-responsiveness. In a second step, we show how to transform V_1 into a VDS V_2 which is secure in the repeated detection game.

A VDS without self-responsiveness. We start by describing a VDS $\mathcal{V}_1 = (\mathtt{Gen}_1, \mathtt{Comp}_1, \mathtt{Chal}_1, \mathtt{Resp}_1, \mathtt{Ver}_1)$ which achieves security without self-responsiveness: The corresponding attack-game $\mathcal{G}^{\mathcal{R},\mathcal{V},W}_{\mathtt{VDS}}$ is derived from the standard attack-game $\mathcal{G}^{\mathcal{R},\mathcal{V},W}_{\mathtt{VDS}}$ by modifying the detection procedure so that instead of emulating \mathcal{R} on the compiled code \widetilde{W} and input (check, c) to compute $y = \mathtt{Resp}(c, \widetilde{W})$, the challenger evaluates $y \stackrel{\$}{\leftarrow} \mathtt{Resp}(c, \widetilde{W})$ himself on the current contents \widetilde{W} of the memory of the emulated RAM.

At a high level, the idea of our construction is as follows: The key generation algorithm \mathtt{Gen}_1 samples a k-bit key K for a symmetric-key cryptosystem, i.e., $K_c = K_v = K \overset{\$}{\leftarrow} \{0,1\}^k$. Given key K, the algorithm \mathtt{Comp}_1 computes an additive sharing $\langle K \rangle$ of K and fills the entire memory \mathtt{MEM} by interleaving a different share of $\langle K \rangle$ between every two words in the original code. Concretely, \mathtt{Comp} compiles some code W for a RAM \mathcal{R} into some new code \widetilde{W} for \mathcal{R} constructed as follows: Between any two consecutive words w_i and w_{i+1} of W the compiler interleaves a uniformly chosen k-bit string $K^{i,i+1} = K_1^{i,i+1} || \dots || K_{\frac{k}{L}}^{i,i+1}$, where each $K_j^{i,i+1} \in \{0,1\}^L$. In the last k bits of the compiled code (i.e., after the last word $w_{|W|-1}$) the string $K^{last} = K \oplus \bigoplus_{i=0}^{|W|-2} K^{i,i+1}$ is written.

To ensure that the compiled code \widetilde{W} executes the same computation as W we need that while being executed it "jumps over" the key-share locations. For this purpose, we do the following modification: After each word w_j of W we insert a (jumpby, n) instruction where n is the number of key-shares between this and the next W-word in the compiled code. Similarly, we modify any "jump" instructions of the original code W to point to the correct locations in \widetilde{W} . Note that this modification is not necessarily applicable to arbitrary code, e.g., for certain self-modifying code it might even be infeasible. However, it is possible for a complete class of code, which we refer to as non-self-modifying structured code, — roughly, this includes code that does not modify its instructions and might only jump by a fixed amount in each operation.

In the following we provide the description of our compiler. The first step to this direction is a process, called Spread, which spreads such a code to allow for enough space between the words to fit the key-shares and adds the extra "jump" instructions to preserve the right program flow. The compiler Comp_1 uses the process Spread to translate, as sketched above, some non-self-modifying structured code W for a RAM $\mathcal{R} = (\mathcal{C}, \operatorname{MEM})$ into (secured) code \widetilde{W} for \mathcal{R} . The algorithm Chal_1 chooses random string $x \in \{0,1\}^k$, and computes the challenge as an encryption of x with key K; i.e., $c \leftarrow \operatorname{Chal}(x,K) = \operatorname{Enc}_K(x)$. For achieving security without self-responsiveness, we can take Enc to be the one-time pad, i.e., $\operatorname{Enc}_K(x) = x \oplus K$. The corresponding algorithm Resp_1 works as follows: On input the challenge $c = \operatorname{Enc}_K(r)$ and the compiled $\operatorname{code} \widetilde{W}$, it reconstructs K by summing up all its shares as retrieved from \widetilde{W} , and outputs a decryption of the challenge under the reconstructed key, i.e., outputs $y = c \oplus K$. The corresponding verification algorithm Ver_1 simply verifies that y = x.

The security of the above scheme follows from the fact that an injection of a sufficiently long continuous virus will have to overwrite a substantial number of bits of a key-share, which will make it infeasible to correctly decrypt the challenge and thus pass the VDS check. The formal security proof of the statement can be found in [23].

¹¹Wlog, here we implicitly assume that $(|\mathtt{MEM}| - |W|) = 0 \mod k$ so that the keys fit exactly in the memory. The general case can also be easily treated.

▶ Theorem 3. Assuming \mathcal{R} is a complete RAM family the VDS \mathcal{V}_1 is secure for \mathcal{R} without self-responsiveness with respect to the class of all non-self-modifying structured codes and the class of adversaries who injects a virus \mathbf{v} with $|\mathbf{v}| \geq 2L + \ell$ where $\ell = \omega(\log k)$ on consecutive memory locations.

Adding Self-Responsiveness. We next modify \mathcal{V}_1 so that it is secure (with self-responsiveness) in the repeated detection game. Due to space limitation, the detailed transformation has been moved to [23]; here we restrict to an overview of the basic technical ideas. The first step is to devise a code W_{Resp} for computing a given VDS response algorithm Resp and then combine it with W via a threading mechanism: When executing, the combined code $\operatorname{Emb}(W,W_{\text{Resp}})$ (periodically) checks if a verification request-input with challenge c is handed to \mathcal{R} ; if so, it stores the CPU state on free memory locations – namely, locations at the end of the memory that are occupied by (no_op) – and changes the program counter to point to the location where W_{Resp} is stored; once Resp has produced output, it restores the CPU to its prior state and continues the execution of W.

But the above modification applied on \mathcal{V}_1 does not yield a secure VDS, as it only detects attacks (virus injections) that occur outside the virus verification procedure, i.e., has the TOCTOU vulnerability discussed in the introduction. To solve this, instead of using a single compilation-key K of length k we use a 2k-bit key K which is parsed as two k-bit keys K_{od} and K_{ev} via the following transformation: Let $K = x_1 || \dots || x_{\frac{2k}{L}}$, where each x_i is a word. Then K_{od} is a concatenation of the odd-indexed words, i.e., x_i 's with $i = 1 \mod 2$, and K_{ev} is a concatenation of the even indexed words. Now the challenge algorithm outputs a double encryption of z with keys K_{od} and K_{ev} , i.e., $c = \text{Enc}_{K_{od}}(\text{Enc}_{K_{ev}}(z))$. In order to decrypt, the response algorithm does the following: First it reconstructs K_{od} by XOR-ing the appropriate shares, and uses it to decrypt c, thus computing $\text{Enc}_{K_{ev}}(z)$. Subsequently, it erases K_{od} from the CPU register (e.g., by filling the register where K_{od} is stored with 0's) and after the erasure completes it starts reconstructing K_{ev} and uses it (as above) to decrypt $\text{Enc}_{K_{ev}}(z)$ and output y = z.

The above modification ensures that in order to correctly answer the challenge, the virus needs both K_{od} and K_{ev} . However, the keys are never simultaneously written in the CPU. Thus if the adversary injects the virus before K_{od} is erased he will overwrite bits from a share of K_{ev} (which at that point exists only in the memory RMEM); thus he will not be able to decrypt the challenge. Otherwise, if the adversary injects the virus after K_{od} has been erased from the CPU, he will overwrite bits from a share of K_{od} ; in this case he might successfully pass this detection attempt, but will fail the next detection attempt.

But we are still not done, because the above argument cannot work with any encryption scheme, e.g., it fails if we instantiate $Enc(\cdot)$ with one-time-pad encryption as in VDS \mathcal{V}_1 . The reason is that once inside the system, the adversary might be able to use the key material it has not destroyed to answer the challenge with non-negligible probability. To avoid this, we make use of a public-key leakage-resilient encryption scheme, e.g., [8], which is secure as long as the adversary's probability of guessing the key is negligible even when one leaks a big part of the key.

The above tricks are the heart of our modified VDS, but a few more low level hacks are employed to ensure that the modified protocols compose well with the ones from the previous section. For example, we need to make sure that the part of the code $\mathrm{Emb}(W,W_{\mathtt{Resp}})$ that

¹² We provide the general solution for k = Lq for some q; with the simplification that L = k, we get that $K = x_1 ||x_2 = K_{od}||K_{ev}$.

32:12 Provably Secure Virus Detection

performs the verification accesses the correct key-share locations. We refer to [23] for the details and the security proof of the resulting VDS $V_2 = (\text{Gen}_2, \text{Comp}_2, \text{Chal}_2, \text{Resp}_2, \text{Ver}_2)$.

On the performance of our VDS. It is easy to verify that the above VDS induces a moderate slowdown, i.e., a factor two, due to processing the newly inserted jump instructions plus whatever slowdown the periodically checking for inputs might impose, on the execution of the code on the RAM (while the code is not being verified). Both slowdowns can be reduced. E.g., the factor two is an overestimate as several instructions are already "jumps" plus reading data does not execute "jumps"; similarly, the check for the input can be executed periodically (or explicitly assumed as being part of any CPU-cycle). On the other hand, verification is quite heavy as it needs to traverse almost the entire memory but this is not our main consideration as (1) we believe that our scheme's provable security guarantee makes the waiting acceptable, and (2) unlike existing attestation methods the verification does not need to stop the normal execution of the program – e.g., in a RAM with parallel processors it could be done by a single designated processor.

In terms of memory usage, the secured code requires a linear (concretely, a factor four) amount of memory with respect to the original code. Although we are less concerned about memory needs – memory is an inexpensive expandable resource the capacities of which increase faster than computation speed – improving the memory usage is an interesting research direction.

5 Detecting Arbitrary Injection

The above VDSs are secure only against (sufficiently) long and consecutive viruses. In this section we describe a construction which is secure independently of the virus length even if the virus bits are not injected continuously, and in particular even when the virus targets non-key-share locations. To achieve this ultimate security guarantees we use a compiler similar to Comp_2 , but we include, for each code word and pair of key-shares, message-authentication-code tags (MACs) which we verify every time we load a code word to the CPU. Concretely, the difference with Comp_2 is that Comp_3 appends a MAC tag t^i to each word, keyed with the (concatenation of the) key-shares $K_{\mathsf{od}}^{i,i+1}||K_{\mathsf{ev}}^{i,i+1}|$ that follow this word. We use a MAC with a special leakage-resilient property that ensures that the adversary cannot forge a tag even when he knows a large portion of the key. Thus, if the malware overwrites only a few bits of some of the key-shares it will be unable to guess an appropriate manipulation of the MAC key. And if it overwrites a large portion it will be unable to answer decryption challenges.

To use the power of the MACs we need to make sure that during the program execution, before loading any word to the CPU we first verify its MAC. To this direction, we assume that, by default, the CPU loads values from the memory RMEM to its registers via a special load instruction (read_auth, i, j), which fetches five consecutive words (corresponding to \widetilde{w}_i and its MAC key and tag), verifies the MAC with the corresponding keys, and only if the MAC verifies, it keeps the word on the CPU to process. If the MAC verification fails, then (read_auth, i, j) deletes at least one of the key-shares from the memory, thus introducing an inconsistency that will be caught by the detection procedure. Furthermore, to ensure that the compiled program will not introduce inconsistencies, our compiler replaces every (write, j, i) instruction (which writes the contents of register j in RMEM location i) with microcode, denoted as write_auth, which, when writing a word in the memory it also updates the corresponding MAC tag. We stress that, unlike read_auth which we need the CPU to support as an atomic instruction, the instruction write_auth does not require any

change in the architecture or additional assumptions as it can be implemented in code itself. The details of the corresponding VDS V_3 along with its security proof can be found in [23].

▶ Theorem 4. Let \mathcal{R} be a complete RAM. If the encryption scheme used in \mathcal{V}_3 is CPA secure even against an adversary who learns all but $\omega(\log k)$ bits of the secret key, then \mathcal{V}_3 is secure for \mathcal{R} in the repeated-detection mode with respect to the class of all non-self-modifying structured codes and adversaries in the bit-by-bit (non-continuous) injection model.

References

- 1 T. AbuHmed, N. Nyamaa, and D. Nyang. Software-based remote code attestation in wireless sensor network. In *GLOBECOM'09*, pages 4680–4687. IEEE Press, 2009.
- W. A. Arbaugh, D. J. Farber, and J. M. Smith. A secure and reliable bootstrap architecture. In SP'97, pages 65-, Washington, DC, USA, 1997. IEEE Computer Society. URL: http://dl.acm.org/citation.cfm?id=882493.884371.
- J. Baron, K. El Defrawy, J. Lampkins, and R. Ostrovsky. How to withstand mobile virus attacks, revisited. In M. M. Halldorsson and S. Dolev, editors, *PODC 2014*, pages 293–302. ACM Press, 2014.
- 4 A. Boldyreva, T. Kim, R. J. Lipton, and B. Warinschi. Provably-secure remote memory attestation for heap overflow protection. Cryptology ePrint Archive, Report 2015/729, 2015
- 5 D. Boneh, R. A. DeMillo, and R. J. Lipton. On the importance of checking cryptographic protocols for faults. In EUROCRYPT' 97, volume 1233 of LNCS, pages 37–51. Spriger, 1997.
- **6** B. Chen and R. Morris. Certifying program execution with secure processors. In *HOTOS'03*, pages 23–23. USENIX Association, 2003.
- 7 D. Dachman-Soled, F.-H. Liu, E. Shi, and H.-S. Zhou. Locally decodable and updatable non-malleable codes and their applications. In *TCC 2015*, volume 9014 of *LNCS*, pages 427–450. Springer, 2015.
- **8** Y. Dodis, S. Goldwasser, Y. T. Kalai, C. Peikert, and V. Vaikuntanathan. Public-key encryption schemes with auxiliary inputs. In *TCC 2010*, volume 5978 of *LNCS*, pages 361–381. Springer, 2010.
- **9** Y. Dodis, Y. T. Kalai, and S. Lovett. On cryptography with auxiliary input. In *STOC '09*, pages 621–630. ACM, 2009.
- 10 S. Dziembowski, K. Pietrzak, and Daniel Wichs. Non-malleable codes. In *ICS 2010*, pages 434–452, 2010.
- 11 K. El Defrawy, G. Tsudik, A. Francillon, and D. Perito. SMART: secure and minimal architecture for (establishing dynamic) root of trust. In NDSS 2012. The Internet Society, 2012.
- 12 C. C. Elgot and A. Robinson. Random-access stored-program machines, an approach to programming languages. *J. ACM*, 11(4):365–399, October 1964.
- S. Faust, P. Mukherjee, J. B. Nielsen, and D. Venturi. A tamper and leakage resilient von neumann architecture. In J. Katz, editor, PKC 2015, volume 9020 of LNCS, pages 579–603. Springer, 2015.
- A. Francillon, Q. Nguyen, K. B. Rasmussen, and G. Tsudik. A minimalist approach to remote attestation. In *DATE'14*, pages 244:1–244:6. European Design and Automation Association, 2014.
- S. Garg, C. Gentry, S. Halevi, M. Raykova, A. Sahai, and B. Waters. Candidate indistinguishability obfuscation and functional encryption for all circuits. In FOCS 2013, pages 40–49. IEEE Computer Society, 2013.

32:14 Provably Secure Virus Detection

- P. Gemmell and M. Naor. Codes for interactive authentication. In *CRYPTO'94*, volume 773 of *LNCS*, pages 355–367. Springer, 1994.
- 17 M. Jakobsson and K.-A. Johansson. Practical and secure Software-Based attestation. In *LightSec 2011*, 2011.
- 18 M. Jakobsson and G. Stewart. Mobile malware: Why the traditional AV paradigm is doomed, and how to use physics to detect undesirable routines. In *BlackHat 2013*, 2013.
- 19 A. Juels and B. S. Kaliski Jr. Pors: Proofs of retrievability for large files. In CCS 2007, pages 584–597. ACM, 2007.
- 20 R. Kennell and L. H. Jamieson. Establishing the genuinity of remote computer systems. In SSYM'03, pages 21–21. USENIX Association, 2003.
- X. Kovah, C. Kallenberg, C. Weathers, A. Herzog, M. Albin, and J. Butterworth. New results for timing-based attestation. In SP 2012, pages 239–253. IEEE Computer Society, 2012.
- 22 R. J. Lipton, R. Ostrovsky, and V. Zikas. Provably secure virus detection. U.S. Patent (pending), Application No. 62/054,160, 2014.
- 23 R. J. Lipton, R. Ostrovsky, and V. Zikas. Provably secure virus detection: Using the observer effect against malware. Cryptology ePrint Archive, Report 2015/728, 2015.
- 24 F. McKeen, I. Alexandrovich, A. Berenzon, C. V. Rozas, H. Shafi, V. Shanbhogue, and U. R. Savagaonkar. Innovative instructions and software model for isolated execution. In HASP'13, pages 10:1–10:1. ACM, 2013.
- W. S. McPhee. Operating systems integrity in os/vs2. *IBM Systems Journal*, 13 Issue 3, pages 230–252, 1974.
- 26 R. Ostrovsky and M. Yung. How to withstand mobile virus attacks (extended abstract). In L. Logrippo, editor, *PODC '91*, pages 51–59. ACM, 1991.
- A. Sahai and B. Waters. How to use indistinguishability obfuscation: deniable encryption, and more. In STOC 2014, pages 475–484. ACM, 2014.
- 28 A. Seshadri, M. Luk, A. Perrig, L. van Doorn, and P. Khosla. Scuba: Secure code update by attestation in sensor networks. In *WiSe'06*, pages 85–94. ACM, 2006.
- 29 A. Seshadri, M. Luk, E. Shi, A. Perrig, L. van Doorn, and P. Khosla. Pioneer: Verifying code integrity and enforcing untampered code execution on legacy systems. In SOSP'05, pages 1–16. ACM, 2005.
- **30** A. Shamir. How to share a secret. Communications of the ACM, 22(11):612–613, 1979.