

Chapter 9

Practical Software Diversification Using In-Place Code Randomization

Vasilis Pappas, Michalis Polychronakis, and Angelos D. Keromytis

Abstract The wide adoption of non-executable page protections has given rise to attacks that employ return-oriented programming (ROP) to achieve arbitrary code execution without the injection of any code. Existing defenses against ROP exploits either require source code or symbolic debugging information, or impose a significant runtime overhead, which limits their applicability for the protection of third-party applications. Aiming for a practical mitigation against ROP attacks, we introduce *in-place code randomization*, a software diversification technique that can be applied directly on third-party software. Our method uses various narrow-scope code transformations that can be applied statically, without changing the location of basic blocks, allowing the safe randomization of stripped binaries even with partial disassembly coverage. We demonstrate how in-place code randomization can prevent the exploitation of vulnerable Windows 7 applications, including Adobe Reader, as well as the automated construction of reliable ROP payloads.

9.1 Introduction

Attack prevention technologies based on the No eXecute (NX) memory page protection bit, which prevent the execution of malicious code that has been injected into a process, are now supported by most recent CPUs and operating systems [49]. The wide adoption of these protection mechanisms has given rise to a new exploitation technique, widely known as *return-oriented programming* (ROP) [61], which allows an attacker to circumvent non-executable page protections without injecting any code. Using return-oriented programming, the attacker can link together small fragments of code, known as *gadgets*, that already exist in the process image of the vulnerable application. Each gadget ends with an indirect control transfer instruction, which transfers control to the next gadget according to a sequence of gadget

Network Security Lab, Columbia University
e-mail: {vpappas, mikepo, angelos}@cs.columbia.edu

addresses injected on the stack or some other memory area. In essence, instead of injecting binary code, the attacker injects just data, which include the addresses of the gadgets to be executed, along with any required data arguments.

Several research works have demonstrated the great potential of return-oriented programming for bypassing defenses such as read-only memory [21], kernel code integrity protections [40], and non-executable memory implementations in mobile devices [29] and operating systems [72, 67, 71, 70]. Consequently, it was only a matter of time for ROP to be employed in real-world attacks. Recent exploits against popular applications, such as the ubiquitous Adobe Reader for Windows, use ROP code to bypass exploit mitigations that are enabled even in the latest OS versions, including Windows 7 SP1. ROP exploits are included in the most common exploit packs, and are actively used in the wild for mounting drive-by download attacks [14, 56].

9.1.1 Existing Defenses

Attackers are able to a priori pick the right code pieces before launching a ROP attack because parts of the code image of the vulnerable application remain static across different installations. Address space layout randomization (ASLR) [49] is meant to prevent this kind of code reuse by randomizing the locations of the executable segments of a running process. However, in both Linux and Windows, parts of the address space do not change due to executables with fixed load addresses [34], or shared libraries incompatible with ASLR [72]. Furthermore, in some exploits, the base address of a DLL can be either calculated dynamically through a leaked pointer [46, 70, 60], or brute-forced [62].

Other defenses against code-reuse attacks complementary to ASLR include compiler extensions [47, 54], code randomization [33, 16, 43], control-flow integrity [11], and runtime solutions [27, 22, 26]. In practice, though, most of these approaches are almost never applied for the protection of the COTS software currently targeted by ROP attacks, either due to the lack of source code or debugging information, or due to their increased overhead. In particular, from the above techniques, those that operate directly on compiled binaries, e.g., by permuting the order of functions [16, 43] or through binary instrumentation [11], require precise and complete extraction of all code and data in the executable sections of the binary. This is possible only if the corresponding symbolic debugging information is available, which however is typically stripped from production binaries.

On the other hand, techniques that do work on stripped binary executables using dynamic binary instrumentation [27, 22, 26], incur a significant runtime overhead that limits their adoption. These defenses are based on monitoring either the frequency of `ret` instructions [22, 26], or the integrity of the stack [27]. At the same time, instruction set randomization (ISR) [42, 13] cannot prevent code-reuse attacks, and current implementations also rely on heavyweight runtime instrumentation or code emulation frameworks.

9.1.2 *In-Place Code Randomization*

Starting with the goal of a practical mitigation against the recent spate of ROP attacks, in this paper we present a novel code randomization method that can harden third-party applications against return-oriented programming. Our approach is based on narrow-scope modifications in the code segments of executables using an array of code transformation techniques, to which we collectively refer as *in-place code randomization* [55]. These transformations are applied statically, in a conservative manner, and modify only the code that can be safely extracted from compiled binaries, without relying on symbolic debugging information. By preserving the length of instructions and basic blocks, these modifications do not break the semantics of the code, and enable the randomization of stripped binaries even without complete disassembly coverage.

The goal of this diversification process is to eliminate or probabilistically modify as many of the gadgets that are available in the address space of a vulnerable process as possible. Since ROP code relies on the correct execution of all chained gadgets, altering the outcome of even a few of them will likely render the ROP code ineffective. The introduced uncertainty raises the bar for the construction of reliable ROP code, as attackers cannot safely assume that a given gadget will perform the intended computation. By randomly choosing and applying different transformations in each instance of the protected application, an attacker will not always be able to choose a safe subset of gadgets that will always remain unchanged.

Still, although quite effective as a standalone mitigation, in-place code randomization is not meant to be a complete prevention solution against ROP exploits, as it offers probabilistic protection and thus cannot deliver any protection guarantees. However, it can be applied in tandem with existing randomization techniques to increase process diversification. This is facilitated by the practically zero overhead of the applied transformations, and the ease with which they can be applied on existing third-party executables.

In the rest of this chapter, we provide some background information on return-oriented programming, discuss the principles of in-place code randomization and the code transformations on which it is based, and present some experimental results using publicly available ROP exploits and automated ROP code generation toolkits.

9.2 From Return-to-Libc to Return-Oriented Programming

9.2.1 *Code-Reuse Attacks*

The introduction of non-executable memory page protections in popular OSes, even for CPUs that do not support the No eXecute (NX) bit, led to the development of the return-to-libc exploitation technique [28]. Using this method, a memory corruption vulnerability can be exploited by transferring control to code that already exists in

the address space of the vulnerable process. By jumping to the beginning of a library function such as `system()`, the attacker can for example spawn a shell without the need to inject any code.

Frequently though, especially for remote exploitation, calling a single function is not enough. In these cases, multiple return-to-libc calls can be “chained” together by ensuring that before returning from one function to the next one, the stack pointer has been correctly adjusted to the beginning of the prepared stack frame for the next call. For instance, for a function with two arguments, this can be achieved by first returning to a short instruction sequence such as `pop reg; pop reg; ret;` found anywhere within the executable part of the process image [53, 52]. The `pop` instructions adjust the stack pointer beyond the arguments of the previously executed function (one `pop` for each argument), and then `ret` transfers control to the next chained function. This approach, however, is not applicable in cases where the function arguments need to be passed through registers. In that case, a few short instruction sequences ending with a `ret` instruction can be chained directly to set the proper registers with the desired arguments, before calling the library function [44].

9.2.2 Return-Oriented Programming

In the above code-reuse techniques, the executed code consists of one or a few short instruction sequences followed by a large block of code belonging to a library function. Hovav Shacham demonstrated that using only a carefully selected set of short instruction sequences ending with a `ret` instruction, known as *gadgets*, it is possible to achieve arbitrary computation, obviating the need for calling library functions [61]. This powerful technique, dubbed *return-oriented programming*, in essence gives the attacker the same level of flexibility offered by arbitrary code injection without injecting any code at all—the injected payload comprises just a sequence of gadget addresses intermixed with any necessary data arguments.

In a typical ROP exploit, the attacker needs to control both the program counter and the stack pointer: the former for executing the first gadget, and the latter for allowing its `ret` instruction to transfer control to subsequent gadgets. Depending on the vulnerability, if the ROP payload is injected in a memory area other than the stack, e.g., the heap, then the stack pointer must first be adjusted to the beginning of the payload through a stack pivot [31, 72]. In a follow up work [20], Checkoway et al. demonstrated that the gadgets used in a ROP exploit need not necessarily end with a `ret` instruction, but with any other indirect control transfer instruction. This also allows the use of any general purpose register in place of the stack pointer as an “index” register for controlling the execution of the gadgets, bypassing any protections based on stack integrity.

Almost a decade after the introduction of the return-to-libc technique [28], the wide adoption of NX-based exploit mitigations in popular OSes sparked a new interest in more advanced forms of code-reuse attacks. The introduction of return-oriented programming [61] and its advancements [19, 21, 40, 20, 29, 17, 59, 66, 72,

71] led to its adoption in real-world attacks [14, 56]. ROP exploits are facilitated by the lack of complete address space layout randomization in both Linux [34], and Windows [72, 41], which otherwise would prevent or at least hinder [62] these attacks.

The ROP code implementations used in recent exploits against Windows applications is mostly based on gadgets ending with `ret` instructions, which conveniently manipulate both the program counter and the stack pointer, although a couple of gadgets ending with `call` or `jmp` are also used for calling library functions. In all publicly available Windows exploits so far, attackers do not have to rely on a fully ROP-based implementation for the whole malicious code that needs to be executed after triggering a memory corruption vulnerability. Instead, ROP code is used only as a first stage for bypassing DEP [49]. Typically, once control flow has been hijacked, the ROP code allocates a memory area with write and execute permissions by calling a library function like `VirtualAlloc`, copies into it some plain shellcode included in the attack vector, and finally jumps to the copied shellcode which now has execute permission [31].

9.3 Approach

In-place code randomization is based on the randomization of the code sections of binary executable files (both libraries and executables) that are part of third-party applications, using an array of binary code transformation techniques. The objective of this randomization process is to break the code semantics of the gadgets that are present in the executable memory segments of a running process, without affecting the semantics of the actual program code.

The execution of a gadget has a certain set of consequences to the CPU and memory state of the exploited process. The attacker chooses how to link the different gadgets together based on which registers, flags, or memory locations each gadget modifies, and in what way. Consequently, the execution of a subsequent gadget depends on the outcome of all previously executed gadgets. Even if the execution of a single gadget has a different outcome than the one anticipated by the attacker, then this will affect the execution of all subsequent gadgets, and it is likely that the logic of the malicious return-oriented code will be severely impacted.

9.3.1 Why In-Place?

The concept of software diversification [23] is the basis for a wide range of protections against the exploitation of memory corruption vulnerabilities. Besides address space layout randomization [49], many techniques focus on the internal randomization of the code segments of executable, and can be combined with ASLR to increase process diversity [33]. Metamorphic transformations [68]. such as the inter-

spersion of ineffectual instructions throughout the code, can shift gadgets from their original offsets and alter many of their instructions, rendering them unusable. Another simpler and probably more effective approach is to rearrange existing blocks of code either at the function level [15, 16, 43, 5], or with finer granularity, at the basic block level [7, 6]. If all blocks of code are reordered so that no one resides at its original location, then all the offsets of the gadgets that the attacker would assume to be present in the code sections of the process will now correspond to completely different code.

These transformations require a precise view of all the code and data objects contained in the executable sections of a PE file, including their cross-references, as existing code needs to be shifted or moved. Due to computed jumps and intermixed data [45], complete disassembly coverage is possible only if the binary contains relocation and symbolic debugging information (e.g., PDB files) [65, 43, 58]. Unfortunately, debugging information is typically stripped from release builds for compactness and intellectual property protection.

For Windows software, in particular, PE files (both DLL and EXE) usually do retain relocation information even if no debugging information has been retained [63]. The loader needs this information in case a DLL must be loaded at an address other than its preferred base address, e.g., because another library has already been mapped to that location, or for ASLR. In contrast to Linux shared libraries and PIC executables, which contain position-independent code and can be easily loaded in an arbitrary location within a process' address space, Windows binaries contain absolute addresses, e.g., as immediate instruction operands or initialized data pointers, that are valid only if the executable has been loaded at its preferred base address. The `.reloc` section of PE files contains a list of offsets relatively to each PE section that correspond to all absolute addresses at which a delta value needs to be added in case the actual load address is different [57].

Relocation information *alone*, however, does not suffice for extracting a complete view of the code within the executable sections of a PE file [7, 65]. Without the symbolic debugging information contained in PDB files, although the location of objects that are reached *only* via indirect jumps *can* be extracted from relocation information, their actual type—code or data—still remains unknown. In some cases, the actual type of these objects could be inferred using heuristics based on constant propagation, but such methods are usually prone to misidentifications of data as code and vice versa. Even a slight shift or size increase of a single object within a PE section will incur cascading shifts to its following objects. Typically, an unidentified object that actually contains code will include PC-relative branches to other code objects. In the absence of the debugging information contained in PDB files, moving such an unidentified code block (or any of its relatively referenced objects) without fixing the immediate displacement operands of all its relative branch instructions that reference other objects, will result to incorrect code.

Given the above constraints, we choose to use only binary code transformations that do not alter the size and location of code and data objects within the executable, allowing the randomization of third-party PE files *without* symbolic debugging information. Although this restriction does not allow us to apply extensive code trans-

formations like basic block reordering or metamorphism, we can still achieve partial code randomization using narrow-scope modifications that can be *safely* applied even without complete disassembly coverage. This can be achieved through slight, in-place code modifications to the correctly identified parts of the code, that do not change the overall structure of basic blocks or functions, but which are enough to alter the outcome of short instruction sequences that can be used as gadgets.

9.3.2 Code Extraction and Modification

Although completely accurate disassembly of stripped x86 binaries is not possible, state-of-the-art disassemblers achieve decent coverage for code generated by the most commonly used compilers, using a combination of different disassembly algorithms [45], the identification of specific code constructs [35], and simple data flow analysis [36]. For our prototype implementation, we use IDA Pro [38] to extract the code and identify the functions of PE executables. IDA Pro is effective in the identification of function boundaries, even for functions with non-contiguous code and extensive use of basic block sharing [39], and also takes advantage of the relocation information present in Windows DLLs.

Typically, however, without the symbolic information of PDB files, a fraction of the functions in a PE executable are not identified, and parts of code remain undiscovered. Our code transformations are applied conservatively, only on parts of the code for which we can be confident that have been accurately disassembled. For instance, IDA Pro speculatively disassembles code blocks that are reached only through computed jumps, taking advantage of the relocation information contained in PE files. However, we do not enable such heuristic code extraction methods in order to avoid any disastrous modifications due to potentially misidentified code. In practice, for the code generated by most compilers, relocation information also ensures that the correctly identified basic blocks have no entry point other than their first instruction. Similarly, some transformations that rely on the proper identification of functions are applied only on the code of correctly recognized functions. Our implementation is separate from the actual code extraction framework used, which means that IDA Pro can be replaced or assisted by alternative code extraction approaches [51, 65, 37], providing better disassembly coverage.

After the code extraction phase is complete, disassembled instructions are first converted to our own internal representation, which holds additional information such as any implicitly used registers, and the registers and flags read or written by the instruction. For correctness, we also track the use of general purpose registers even in floating point, MMX, and SSE instructions. Although these type of instructions have their own set of registers, they do use general purpose registers for memory references (e.g., as the `fmul` instruction in Fig. 9.1). We then proceed and apply the in-place code transformations discussed in the following section. These are applied only on the parts of the executable segments that contain (intended or unintended [61]) instruction sequences that can be used as gadgets. As a result of

some of the transformations, instructions may be moved from their original locations within the same basic block. In these cases, for instructions that contain an absolute address in some of their operands, the corresponding entries in the `.reloc` sections of the randomized PE file are updated with the new offsets where these absolute addresses are now located.

9.3.3 Deployment

Our publicly-available prototype implementation, called `orp`, can be used to generate randomized instances of existing applications. `Orp` processes each PE file individually, and generates multiple randomized copies that can then replace the original. To relieve users of the burden of installing and running `orp`, as part of our future work we also plan to create a web service that will allow the submission of executables for randomization.

Given the complexity of the analysis required for generating a set of randomized instances of an input file (in the order of a few minutes on average for the PEs used in our tests), `orp` can be used for the off-line generation of a pool of randomized PE files for a given application. Note that for most of the tested Windows applications, only some of the DLLs need to be randomized, as the rest are usually ASLR-enabled (although they can also be randomized for increased protection). In a production deployment, a system service or a modified loader can then pick a different randomized version of the required DLLs and executables each time the application is launched, following the same way of operation as tools like EMET [48].

9.4 In-Place Code Transformations

In this section we present in detail the different code transformations used for in-place code randomization. Although some of the transformations such as instruction reordering and register reassignment are also used by compilers and polymorphic code engines for code optimization [12] and obfuscation [68], applying them at the binary level—without having access to the higher-level structural and semantic information available in these settings—poses significant challenges.

9.4.1 Atomic Instruction Substitution

One of the basic concepts of code obfuscation and metamorphism [68] is that the exact same computation can be achieved using a countless number of different instruction combinations. When applied for code randomization, substituting the instructions of a gadget with a functionally-equivalent—but different—sequence of

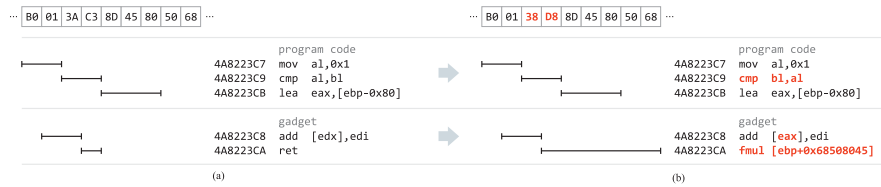


Fig. 9.1: Example of atomic instruction substitution. The equivalent, but different form of the `cmp` instruction does not change the original program code (a), but renders the non-intended gadget unusable (b).

instructions would not affect any ROP code that uses that gadget, since its outcome would be the same. However, by modifying the instructions of the original program code, this transformation in essence modifies certain bytes in the code image of the program, and consequently, can drastically alter the structure of non-intended instruction sequences that overlap with the substituted instructions.

Many of the gadgets used for return-oriented programming consist of unaligned instructions that have not been emitted by the compiler, but which happen to be present in the code image of the process due to the density and variable-length nature of the x86 instruction set. In the example of Fig. 9.1(a), the actual code generated by the compiler consists of the instructions `mov`; `cmp`; `lea`; starting at byte B0.¹ However, when disassembling from the next byte, a useful non-intended gadget ending with `ret` is found.

Compiled code is highly optimized, and thus the replacement of even a single instruction in the original program code usually requires either a longer instruction, or a combination of more than one instruction, for achieving the same purpose. Given that our aim is to randomize the code of stripped binaries, even a slight increase in the size of a basic block is not possible, which makes the most commonly used instruction substitution techniques unsuitable for our purpose.

In certain cases though, it is possible to replace an instruction with a single, functionally-equivalent instruction of the *same* length, thanks to the flexibility offered by the extensive x86 instruction set. Besides obvious candidates based on replacing addition with negative subtraction and inversely, there are also some instructions that come in different forms, with different opcodes, depending on the supported operand types. For example, `add r/m32, r32` stores the result of the addition in a register *or* memory operand (*r/m32*), while `add r32, r/m32` stores the result in a register (*r32*). Although these two forms have different opcodes, the two instructions are equivalent when both operands happen to be registers. Many arithmetic and logical instructions have such dual equivalent forms, while in some cases there can be up to five equivalent instructions (e.g., `test r/m8, r8`, `or r/m8, r8`, `or r8, r/m8`, and `r/m8, r8`, and `r8, r/m8`, affect the flags of the EFLAGS

¹ The code of all examples throughout this chapter comes from `icucnv36.dll`, included in Adobe Reader v9.3.4. This DLL was used for the ROP code of a DEP-bypass exploit for CVE-2010-2883 [1] (see Table 9.2).

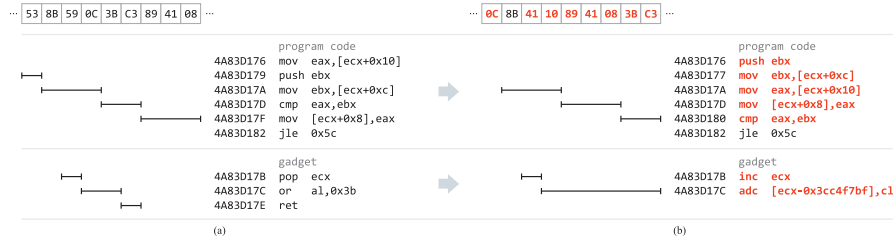


Fig. 9.2: Example of how intra basic block instruction reordering can affect a non-intended gadget.

register in the same way when both operands are the *same* register). In our prototype implementation we use the sets of equivalent instructions used in Hydan [30], a tool for hiding information in x86 executables, with the addition of one more set that includes the equivalent versions of the `xchg` instruction.

As shown in Fig. 9.1(b), both operands of the `cmp` instruction are registers, and thus it can be replaced by its equivalent form, which has different opcode and Mod-R/M bytes [10]. Although the actual program code does not change, the `ret` instruction that was “included” in the original `cmp` instruction has now disappeared, rendering the gadget unusable. In this case, the transformation completely *eliminates* the gadget, and thus will be applied in all instances of the randomized binary. In contrast, when a substitution does not affect the gadget’s final indirect jump, then it is applied probabilistically.

9.4.2 Instruction Reordering

In certain cases, it is possible to reorder the instructions of small self-contained code fragments without affecting the correct operation of the program. This transformation can significantly impact the structure of non-intended gadgets, but can also break the attacker’s assumptions about gadgets that are part of the actual code.

9.4.2.1 Intra Basic Block Reordering

The actual instruction scheduling chosen during the code generation phase of a compiler depends on many factors, including the cost of instructions in cycles, and the applied code optimization techniques [12]. Consequently, the code of a basic block is often just one among several possible instruction orderings that are all equivalent in terms of correctness. Based on this observation, we can partially modify the code within a basic block by reordering some of its instructions according to an alternative instruction scheduling.

The basis for deriving an alternative instruction scheduling is to determine the ordering relationships among the instructions, which must always be satisfied to maintain code correctness. The *dependence graph* of a basic block represents the instruction interdependencies that constrain the possible instruction schedules [50]. Since a basic block contains straight-line code, its dependence graph is a directed acyclic graph with machine instructions as vertices, and dependencies between instructions as edges. We apply dependence analysis on the code of disassembled basic blocks to build their dependence graph using an adaptation of a standard dependence DAG construction algorithm [50, Fig. 9.6] for machine code. Applying dependence analysis directly on machine code requires a careful treatment of the dependencies between x86 instructions. Compared to the analysis of code expressed in an intermediate representation form, this includes the identification of data dependencies not only between register and memory operands, but also between CPU flags and implicitly used registers and memory locations.

For each instruction i , we derive the sets $use[i]$ and $def[i]$ with the registers used and defined by the instruction. Besides register operands and registers used as part of effective address computations, this includes any implicitly used registers. For example, the use and def sets for `pop eax` are $\{esp\}$ and $\{eax, esp\}$, while for `rep stosb`² are $\{ecx, eax, edi\}$ and $\{ecx, edi\}$, respectively. We initially assume that all instructions in the basic block depend on each other, and then check each pair for read-after-write (RAW), write-after-read (WAR), and write-after-write (WAW) dependencies. For example, i_1 and i_2 have a RAW dependency if any of the following is true: i) $def[i_1] \cap use[i_2] \neq \emptyset$, ii) the destination operand of i_1 and the source operand of i_2 are both a memory location, iii) i_1 writes at least one flag read by i_2 .

Note that condition ii) is quite conservative, given that i_2 will actually depend on i_1 only if i_2 reads the *same* memory location written by i_1 . However, unless both memory operands use absolute addresses, it is hard to determine statically if the two effective addresses point to the same memory location. In our future work, we plan to use simple data flow analysis to relax this condition. Besides instructions with memory operands, this condition should also be checked for instructions with implicitly accessed memory locations, e.g., `push` and `pop`. The conditions for WAR and WAW dependencies are analogous. If no conflict is found between two instructions, then there is no constraint in their execution order.

Figure 9.2(a) shows the code of a basic block that contains a non-intended gadget, and Fig. 9.3 its corresponding dependence DAG. Instructions not connected via a direct edge are independent, and have no constraint in their relative execution order. Given the dependence DAG of a basic block, the possible orderings of its instructions correspond to the different topological sorting arrangements of the graph [69]. Fig. 9.2(b) shows one of the possible alternative orderings of the original code. The locations of all but one of the instructions and the values of all but one of the bytes have changed, eliminating the non-intended gadget contained in the original code. Although a new gadget has appeared a few bytes further into the block, (ending

² `stosb` (Store Byte to String) copies the least significant byte from the `eax` register to the memory location pointed by the `edi` register and increments `edi`'s value by one. The `rep` prefix repeats this instruction until `ecx`'s value reaches zero, while decreasing it after each repetition.

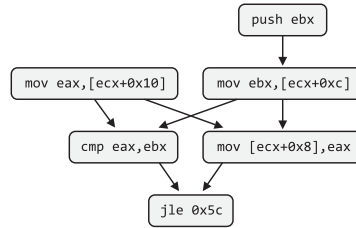


Fig. 9.3: Dependence graph of the basic block shown in Fig. 9.2.

again with a `ret` instruction at byte C3), an attacker cannot depend on it since alternative orderings will shift it to other locations, and some of its internal instructions will always change (e.g., in this example, the useful `pop ecx` is gone). In fact, the `ret` instruction can be eliminated altogether using atomic instruction substitution.

An underlying assumption we make here is that basic block boundaries will not change at runtime. If a computed control transfer instruction targets a basic block instruction other than its first, then reordering may break the semantics of the code. Although this may seem restrictive, we note that throughout our evaluation we did not encounter any such case. For compiler-generated code, IDA Pro is able to compute all jump targets even for computed jumps based on the PE relocation information. In the most conservative case, users may choose to disable instruction reordering and still benefit from the randomization of the other techniques—Section 9.5 includes results for each technique individually.

9.4.2.2 Reordering of Register Preservation Code

The calling convention followed by the majority of compilers for Windows on x86 architectures, similarly to Linux, specifies that the `ebx`, `esi`, `edi`, and `ebp` registers are callee-saved [32]. The remaining general purpose registers, known as scratch or volatile registers, are free for use by the callee without restrictions. Typically, a function that needs to use more than the available scratch registers, preserves any non-volatile registers before modifying them by storing their values on the stack. This is usually done at the function prologue through a series of `push` instructions, as in the example of Fig. 9.4(a), which shows the very first and last instructions of a function. At the function epilogue, a corresponding series of `pop` instructions restores the saved values from the stack, right before returning to the caller.

Sequences that contain `pop` instructions followed by `ret` are among the most widely used gadgets found in ROP exploits, since they allow the attacker to load registers with values that are supplied as part of the injected payload [64]. The order of the `pop` instructions is crucial for initializing each register with the appropriate value. For example, loading `01020304` to `esi` and `DEADCODE` to `ebx` using the gadget `pop esi; pop ebx; ret;` found in the epilogue of the function in Fig. 9.4, would require the following arrangement in the ROP payload:

<pre> 4A834B3B 0 push ebx 4A834B3C -4 push esi 4A834B3D -8 mov ebx,ecx 4A834B3F -8 push edi 4A834B40 -C mov esi,edx ... 4A834B7C -C pop edi 4A834B7D -8 pop esi 4A834B7E -4 pop ebx 4A834B7F 0 ret </pre> <p style="text-align: center;">(a)</p>	➔	<pre> 4A834B3B push edi 4A834B3C push ebx 4A834B3D push esi 4A834B3E mov ebx,ecx 4A834B40 mov esi,edx ... 4A834B7C pop esi 4A834B7D pop ebx 4A834B7E pop edi 4A834B7F ret </pre> <p style="text-align: center;">(b)</p>
--	---	---

Fig. 9.4: Example of instruction reordering in the register preservation code at the preamble and epilogue of a function.

```

.. |7D 6B 83 4A|04 03 02 01|DE C0 AD DE|B3 02 83 4A| ..
| gdtg addr |   esi   |   ebx   | next gdtg |

```

As seen in the function prologue, the compiler stores the values of the callee-saved registers in arbitrary order, and sometimes the relevant `push` instructions are interleaved with instructions that use previously-preserved registers. At the function epilogue, the saved values are `pop`'ed from the stack in reverse order, so that they end up to the proper register. Consequently, as long as the saved values are restored in the right order, their actual order on the stack is irrelevant. Based on this observation, we can randomize the order of the `push` and `pop` instructions of register preservation code by maintaining the first-in-last-out order of the stored values, as shown in Fig. 9.4(b). In this example, there are six possible orderings of the three `pop` instructions, which means that any assumption that the attacker may make about which registers will hold the two supplied values, will be correct with a probability of one in six (or one in three, if only one register needs to be initialized). In case only two registers are preserved, there are two possible orderings, allowing the gadget to operate correctly half of the time.

This transformation is applied conservatively, only to functions with accurately disassembled prologue and epilogue code. To make sure that we properly match the `push` and `pop` instructions that preserve a given register, we monitor the stack pointer delta throughout the whole function, as shown in the second column of Fig. 9.4(a). If the deltas at the prologue and epilogue do not match, e.g., due to call sites with unknown calling conventions throughout the function, or indirect manipulation of the stack pointer, then no randomization is applied. As shown in Fig. 9.4(b), any non-preservation instructions in the function prologue are reordered along with the `push` instructions by maintaining any interdependencies, as discussed in the previous section. For functions with multiple exit points, the preservation code at all epilogs should match the function's prologue. Note that there can be multiple `push` and `pop` pairs for the same register, in case the register is preserved only throughout some of the execution paths of a function.

9.4.3 Register Reassignment

During the register allocation phase, the compiler assigns the arbitrarily many variables of the higher-level program into the much smaller set of registers that are available in the target processor architecture. Although the program points at which a certain variable should be stored in a register or spilled into memory are chosen by sophisticated allocation algorithms, the actual name of the general purpose register that will hold a particular variable is mostly an arbitrary choice. That is, whenever a new variable needs to be mapped to a register, the compiler can pick any of the available registers at that point to hold it. As a result, the actual register assignment throughout the code of a given compiled binary is just one among many possible register assignments. Based on this observation, we can reassign the names of the register operands in the existing code according to a different—but equivalent—register assignment, without affecting the semantics of the original code. When considering each gadget as an autonomous code sequence, this transformation can alter the outcome of many gadgets, which will now read or modify different registers than those assumed by the attacker.

Due to the much higher cost of memory accesses compared to register accesses, compilers strive to map as many variables as possible to the available registers. Consequently, at any point in a large program, multiple registers are usually in use, or *live* at the same time. Given the control flow graph (CFG) of a compiled program, a register r is *live* at a program point p iff there is a path from p to a use of r that does not go through a definition of r . The *live range* of r is defined as the set of program points where r is live, and can be represented as a subgraph of the CFG [18]. Since the same register can hold different variables at different points in the program, a register can have multiple disjoint live regions in the same CFG.

For each correctly identified function, we compute the live ranges of all registers used in its body by performing liveness analysis [12] directly on the machine code. Given the CFG of the function and the sets $use[i]$ and $def[i]$ for each instruction i , we derive the sets $in[i]$ and $out[i]$ with the registers that are *live-in* and *live-out* at each instruction. For this purpose, we use a modified version of a standard live-variable analysis algorithm [12, Fig. 9.16] that computes the *in* and *out* sets at the instruction level, instead of the basic block level. The algorithm computes the two sets by iteratively reaching a fixed point for the following data-flow equations: $in[i] = use[i] \cup (out[i] - def[i])$ and $out[i] = \bigcup \{in[s] : s \in succ[i]\}$, where $succ[i]$ is the set of all possible successors of instruction i .

Figure 9.5 shows part of the CFG of a function and the corresponding live ranges for `eax` and `edi`. Initially, we assume that all registers are live, since some of them may hold values that have been set by the caller. In this example, `edi` is live when entering the function, and the `push` instruction at line 2 stores (uses) its current value on the stack. The following `mov` instruction initializes (defines) `edi`, ending its previous live range (d_0). Note that although a live range is a sub-graph of the CFG, we illustrate and refer to the different live ranges as linear regions for the sake of convenience.

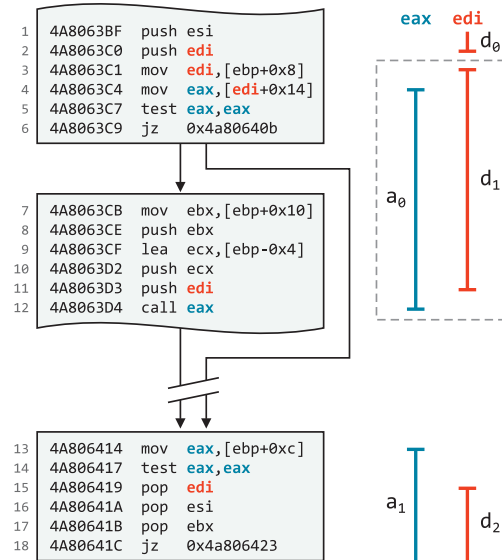


Fig. 9.5: The live ranges of `eax` and `edi` in part of a function. The two registers can be swapped in all instructions throughout their parallel, self-contained regions a_0 and d_1 (lines 3–12).

The next definition of `edi` is at line 15, which means that the last use of its previous value at line 11 also ends its previous live region d_1 . Region d_1 is a *self-contained* region, within which we can be confident that `edi` holds the same variable. The `eax` register also has a self-contained live region (a_0) that runs in parallel with d_1 . Conceptually, the two live ranges can be extended to share the same boundaries. Therefore, the two registers can be swapped across all the instructions located within the boundaries of the two regions, without altering the semantics of the code.

The `call eax` instruction at line 12 can be conveniently used by an attacker for calling a library function or another gadget. By reassigning `eax` and `edi` across their parallel live regions, any ROP code that would depend on `eax` for transferring control to the next piece of code, will now jump to an incorrect memory location, and probably crash. For code fragments with just two parallel live regions, an attacker can guess the right register half of the times. In many cases though, there are three or more general purpose registers with parallel live regions, or other available registers that are live before or after another register’s live region, allowing for a higher number of possible register assignments.

The registers used in the original code can be reassigned by modifying the Mod-R/M and sometimes the SIB byte of the relevant instructions. As in previous code transformations, besides altering the operands of instructions in the existing code, these modifications can also affect overlapping instructions that may be part of non-intended gadgets. Note that implicitly used registers in certain instructions cannot be replaced. For example, the one-byte “move data from string to string” instruction

(`movs`) always uses `esi` and `edi` as its source and destination operands, and there is no other one-byte instruction for achieving the same operation using a different set of registers [10]. Consequently, if such an instruction is part of the live region of one of its implicitly used registers, then this register cannot be reassigned throughout that region. For the same reason, we exclude `esp` from liveness analysis.

Finally, although calling conventions are followed for most of the functions, this is not always the case, as compilers are free to use any custom calling convention for private or static functions. Most of these cases are conservatively covered through a bottom-up call analysis that discovers custom register arguments and return value registers. First, all the external function definitions found in the import table of the DLL are marked as level-0 functions. IDA Pro can effectively distinguish between different calling conventions that these external functions may follow, and reports their declaration in the C language. Thus, in most cases, the register arguments and the return value register (if any) for each of the level-0 functions are known. For any `call` instruction to a level-0 function, its register arguments are added to `call`'s set of implicitly read registers, and its return value registers are added to `call`'s set of implicitly written registers.

In the next phase, level-1 functions are identified as the set of functions that call only level-0 functions or no other function. Any registers read by a level-1 function, without prior writing them, are marked as its register arguments. Similarly, any registers written and not read before a return instruction are marked as return value registers. Again, the sets of implicitly read and written register of all the `call` instructions to level-1 functions are updated accordingly. Similarly, level-2 functions are the ones that call level-1 or level-0 functions, or no other function, and so on. The same process is repeated until no more function levels can be computed. The intuition behind this approach is that private functions, which may use non-standard calling conventions, are called by other functions in the same DLL and, in most cases, not through computed call instructions.

9.5 Randomization Analysis

A crucial aspect for the effectiveness of in-place code randomization is the randomization coverage in terms of what percentage of the gadgets found in an executable can be safely randomized. A gadget may remain intact for one of the following reasons: i) it is part of data embedded in a code segment, ii) it is part of code that could not be disassembled, or iii) it is not affected by any of our transformations. In this section, we explore the randomization coverage of our prototype implementation using a large data set of 5,235 PE files (both DLL and EXE), detailed in Table 9.1.

For each PE file, we first pinpoint all gadgets contained in its executable sections. We consider as a gadget [61] any intended or unintended instruction sequence that ends with an indirect control transfer instruction, and which does not contain i) a privileged or invalid instruction (can occur in non-intended instruction sequences), and ii) a control transfer instruction other than its final one, with the exception of

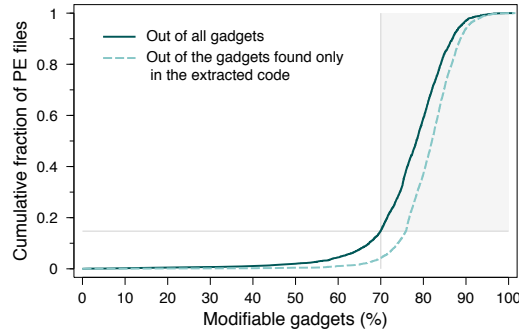


Fig. 9.6: Percentage of modifiable gadgets for a set of 5,235 PE files (detailed in Table 9.1). Indicatively, for the upper 85% of the files, more than 70% of *all* gadgets in the executable segments of each PE file can be modified (shaded area).

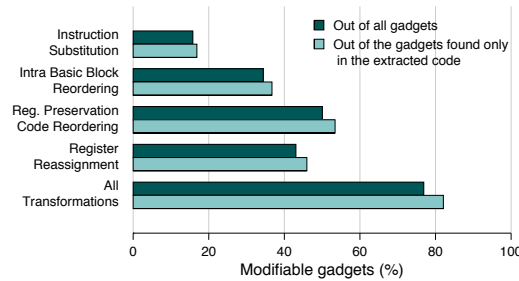


Fig. 9.7: Percentage of modifiable gadgets according to the different code transformations.

Table 9.1: Modifiable (eliminated vs. broken) gadgets for a collection of various PE files.

Software	PE Files	Total	Modifiable %	Eliminated %	Broken %
Adobe Reader 9	43	1,250,959	75.4	8.7	66.7
Firefox 4	28	458,760	83.0	12.4	70.6
iTunes 10	75	396,478	74.0	8.0	66.0
Windows XP SP3	1,698	8,305,177	77.7	9.3	68.4
Windows 7 SP1	3,391	16,951,300	76.5	9.7	66.8
Total	5,235	27,362,674	76.9	9.5	67.4

indirect `call` (can be used in the middle of a gadget for calling a library function). We assume a maximum gadget length of five instructions, which is typical for existing ROP code implementations [61, 20]. The larger the length of the gadget, the higher the probability that at least one of its instructions will be affected. However, for larger gadgets, it is possible that the modified part of the gadget may be irrelevant

for the purpose of the attacker. For example, if only the first instruction of the gadget `inc eax; pop ebx; ret;` is randomized, this will not affect any ROP code that either does not rely on the value of `eax` at that point, or uses the shorter gadget `pop ebx; ret;` directly. For this reason, we consider *all* different subsequences with length between two to five instructions as separate gadgets.

Figure 9.6 shows the percentage of modifiable gadgets out of *all* gadgets found in the executable sections of each PE file (solid line), as a cumulative fraction of all PE files in the data set. In about 85% of the PE files, more than 70% of the gadgets can be randomized by our code transformations. Many of the unmodified gadgets are located in parts of code that have not been extracted by IDA Pro, and which consequently will never be affected by our transformations. When considering only the gadgets that are contained within the disassembled code regions on which code randomization can be applied, the percentage of affected gadgets slightly increases (dashed line). Given that we do not take into account code blocks that have been identified by IDA Pro using speculative methods, this shows that the use of a more sophisticated code extraction mechanism will increase the number of gadgets that can be modified. Figure 9.7 shows the total percentage of gadgets modified by each code transformation technique for the same data set. Note that a gadget can be modified by more than one technique. Overall, the total percentage of modifiable gadgets across all PE files is about 76.9%, as shown in Table 9.1.

We identify two qualitatively different ways in which a code transformation can impact a gadget. As discussed in Sec. 9.4.1, a gadget can be *eliminated*, if any of the applied transformations removes completely its final control transfer instruction. If the final control transfer instruction remains intact, a gadget can then be *broken*, if at least one of its internal instructions is altered, and the CPU and memory state after its execution is different than the original, i.e., the outcome of its computation is not the same. As shown in Table 9.1, in the average case, about 9.5% of *all* gadgets contained in a PE file can be rendered completely unusable. For a vulnerable application, this already removes about one in ten of the available gadgets for the construction of ROP code. Although the rest of the modifiable gadgets (67.4%) is not eliminated, they can be “broken” by probabilistically modifying one or more of their instructions.

9.6 Correctness and Performance

One of the basic principles of our approach is that the different in-place code randomization techniques should be applied cautiously, without breaking the semantics of the program. A straightforward way to verify the correctness of our code transformations is to apply them on existing code and compare the outcome before and after modification. Simply running a randomized version of a third-party application and verifying that it behaves in the expected way can provide a first indication. However, using this approach, it is hard to exercise a significant part of the code, and potentially incorrect modifications may go unnoticed.

Table 9.2: ROP exploits [1, 3, 4] and generic ROP payloads [8, 24] tested on Windows 7 SP1.

ROP exploit/payload	Gadgets in		Modifiable (total %: Broken % Eliminated %)	Unique Gadgets		Combinations
	non-ASLR DLLs			Used: Modifiable (Broken, Elim.)		
Adobe Reader [1]	36,760	28,637 (77.9: 70.1 7.8)	11: 6 (5, 1)		287	
Integard Pro [3]	5,137	4,027 (78.4: 70.5 7.9)	16: 10 (9, 1)		322,559	
Mplayer Lite [4]	117,822	104,671 (88.8: 70.0 18.8)	18: 7 (6, 1)		1,128,959	
msvcr71.dll [8]	10,301	7,129 (69.2: 59.6 9.6)	14: 9 (8, 1)		3,317,760	
msvcr71.dll [24]	10,301	7,129 (69.2: 59.6 9.6)	16: 8 (8, 0)		1,728,000	
mscorie.dll [8]	1,616	1,304 (80.6: 73.5 7.1)	10: 4 (4, 0)		25,200	
mfc71u.dll [24]	86,803	64,053 (73.8: 68.7 5.1)	11: 6 (6, 0)		170,496	

For this purpose, we used the test suite of Wine [9], a compatibility layer that allows Windows applications to run on Unix-like operating systems. Wine provides alternative implementations of the DLLs that comprise the Windows API, and comes with an extensive test suite that covers the implementations of most of the functions exported by the core Windows DLLs. Each function is executed multiple times using various inputs that test different conditions, and the outcome of each execution is compared against a known, expected result. We ported the test code for about one third of the 109 DLLs included in the test suite of Wine v1.2.2, and used it directly on the actual DLLs gathered from a Windows 7 installation. Using multiple randomized versions of each tested DLL, we verified that in all runs, all tests completed successfully.

We took advantage of the extensive and diverse code execution coverage of this experiment to also evaluate the impact of in-place code randomization to the runtime performance of the modified code. Among the different code transformations, instruction reordering is the only one that could potentially introduce some non-negligible overhead, given that sometimes the chosen ordering may be sub-optimal. We measured the overall CPU user time for the completion of all tests by taking the average time across multiple runs, using both the original and the randomized versions of the DLLs. In all cases, there was no observable difference in the two times, within measurement error.

9.7 Effectiveness Against Real-World ROP Exploits

9.7.1 ROP Exploits and Generic ROP Payloads

We evaluated the effectiveness of in-place code randomization using publicly available ROP exploits against vulnerable Windows applications [1, 3, 4], as well as generic ROP payloads based on commonly used DLLs [8, 24]. These seven different ROP code implementations, listed in Table 9.2, bypass Windows DEP and execute a

second-stage shellcode, as described in Sec. 9.2, and work even in the latest version of Windows, with DEP and ASLR enabled. We first verified that all exploits and payloads succeed by testing them against installations of the vulnerable applications on a Windows 7 SP1 virtual machine, using a shellcode that just spawns `calc.exe`. The ROP code used in the three exploits is implemented with gadgets from one or a few DLLs that do not support ASLR, as shown in the second column of Table 9.2. The number of unique gadgets used in each case varies between 10–18, and typically a large part of the gadgets is repeatedly executed at many points throughout the ROP code. When replacing the original non-ASLR DLLs of each application with randomized versions, in all cases the exploits were rendered unsuccessful. Similarly, we used a custom application to test the generic ROP payloads and verified that the ROP code did not succeed when the corresponding DLL was randomized.

The ROP code of the exploit against Acrobat Reader uses just 11 unique gadgets, all coming from a single non-ASLR DLL (`icucnv36.dll`). From these gadgets, in-place code randomization can alter six of them: one gadget is completely eliminated, while the other five broken gadgets have 2, 2, 3, 4, and 6 possible states, respectively, resulting to a total of 287 randomized states (*in addition* to the always eliminated gadget, which also alone breaks the ROP code). Even if we assume that no elimination were possible, the exploit would still succeed only in one out of the 288 (0.35%) possible instances (including the original) of the given gadget set. Considering that this is a client-side exploit, in which the attacker will probably have only one or a few opportunities for tricking the user to open the malicious PDF file, the achieved randomization entropy is quite high—always assuming that none of the gadgets could have been eliminated. As shown in Table 9.2, the number of possible randomized states in the rest of the cases is several orders of magnitude higher. This is mostly due to the larger number of broken gadgets, as well as due to a few broken gadgets with tens of possible modified states, which both increase the number of states exponentially.

Next, we explored whether the affected gadgets could be directly replaced with unmodifiable gadgets in order to reliably circumvent our technique. Out of the six affected gadgets in the Adobe Reader exploit, only four can be directly replaced, meaning that the exploit cannot be trivially modified to bypass randomization. Furthermore, two of the gadgets have only one replacement each, and both replacements are found in code regions that are not discovered by IDA Pro—both could be randomized using a more precise code extraction method. For the rest of the ROP payloads, there are at least three irreplaceable gadgets in each case.

We should note that the relatively small number of gadgets used in most of these ROP payloads is a worst-case scenario for our technique, which however not only is able to prevent these exploits, but also does not allow the attacker to directly replace all the affected gadgets. Indeed, besides the more complex ROP payloads used in the Integard and Mplayer exploits, the rest of the payloads use API functions that are already imported by a non-ASLR DLL, and simply call them directly using hard-coded addresses. This type of API invocation is much simpler and requires fewer gadgets [59] compared to ROP code like the one used in the Integard and Mplayer exploits (16 and 18 unique gadgets, respectively), which first dynamically locates a

pointer to `kernel32.dll` (always ASLR-enabled in Windows 7) and then gets a handle to `VirtualProtect`.

9.7.2 *Hindering Automated ROP Payload Generation*

The fact that some of the randomized gadgets are not directly replaceable does not necessarily mean that the same outcome cannot be achieved using solely unmodifiable gadgets. For example, a gadget that performs an arithmetic operation and then copies the result to a memory location could be trivially replaced with two gadgets: one that does the arithmetic operation and one that copies the result. To assess whether an attacker could construct a ROP payload resistant to in-place code randomization based on gadgets that cannot be randomized, we used Q [59] and Mona [25], two automated ROP code construction tools.

Q is a general-purpose ROP compiler that uses semantic program verification techniques to identify the functionality of gadgets, and provides a custom language, named QooL, for writing input programs. Its current implementation only supports simple QooL programs that call a single function or system call, while passing a single custom argument. In case the function to be called belongs to an ASLR-enabled DLL, Q can compute a handle to it through the import table of a non-ASLR DLL [34], when applicable. We should note that although Q currently compiles only basic QooL programs that call a single API function, this does not limit our evaluation, but on the contrary, stresses even more our technique. The simpler the programs, the fewer the gadgets used, which makes it easier for Q to generate ROP code even when our technique limits the number of available gadgets.

Mona is a plug-in for Immunity Debugger [2] that automates the process of building Windows ROP payloads for bypassing DEP. Given a set of non-ASLR DLLs, Mona searches for available gadgets, categorizes them according to their functionality, and then attempts to automatically generate four alternative ROP payloads for giving execute permission to the embedded shellcode and then invoking it, based on the `VirtualProtect`, `VirtualAlloc`, `NtSetInformationProcess`, and `SetProcessDEPPolicy` API functions (the latter two are not supported in Windows 7).

Considering the functionality of the ROP payloads generated by the two tools, Mona generates slightly more complex payloads, but its gadget composition engine is less sophisticated compared to Q's. Q generates payloads that compute a function address, construct its single argument, and call it. Payloads generated by Mona also call a single memory allocation API function (which though requires the construction of several arguments), copy the shellcode to the newly allocated area, and transfer control to it. Note that the complexity of the ROP code used in the tested exploits is even higher, since they rely on up to four different API functions [1], or "walk up" the stack to discover pointers to non-imported functions from ASLR-enabled DLLs [3, 4].

Table 9.3: Results of running Q [59] and Mona [25] on the original non-ASLR DLLs listed in Table 9.2, and the unmodified parts of their randomized versions. In all cases, both tools failed to generate a ROP payload using solely non-randomized gadgets.

Application/DLL	Q success		Mona success	
	Orig.	Rand.	Orig.	Rand.
Adobe Reader	✓	✗	✓ (VA)	✗
Integard Pro	✓	✗	✗	✗
Mplayer	✓	✗	✓ (VA)	✗
msvcr71.dll	✓	✗	✗	✗
mscorie.dll	✗	✗	✗	✗
mfc71u.dll	✓	✗	✓ (VA, VP)	✗

Table 9.3 shows the results of running Q and Mona on the same set of applications and DLLs used in the previous section (for applications, all non-ASLR DLLs are analyzed collectively), for two different cases: when all gadgets are available to the ROP compiler, and when only the non-randomized gadgets are available. The second case aims to build a payload that will be functional even when code randomization is applied. Although both Q and Mona were able to create payloads when applied on the original DLLs in almost all cases, they failed to construct any payload using only non-randomized gadgets in *all* cases.

Although our technique was able to prevent two different tools from automatically constructing reliable ROP code, this favorable outcome does not exclude the possibility that a functional payload could still be constructed based solely on non-randomized gadgets, e.g., in a manual way or using an even more sophisticated ROP compiler. However, it clearly demonstrates that in-place code randomization significantly raises the bar for attackers, and makes the construction of reliable ROP code much harder, even in an automated way.

9.8 Discussion

Randomization Coverage. In-place code randomization may not always randomize a significant part of the executable address space, and it is hard to give a definitive answer on whether the remaining unmodifiable gadgets, or even some of the partially affected gadgets, would be sufficient for constructing useful ROP code. This depends on the code in the non-ASLR address space of the particular vulnerable process, as well as on the actual operations that need to be achieved using ROP code. Note that Turing-completeness is irrelevant for practical exploitation [59], and none of the gadget sets used in the tested ROP payloads is Turing-complete. For this reason, we emphasize that in-place code randomization should be used as a mitigation technique, in the same fashion as application armoring tools like EMET [48], and not as a complete prevention solution.

As previous studies [59, 61, 29] have shown, though, the feasibility of building a ROP payload is proportional to the size of the non-ASLR code base, and reversely proportional to the complexity of the desired functionality. Our experimental evaluation shows that in all cases, the space of the remaining useful gadgets after randomization is sufficiently small to prevent the automated generation of ROP payloads. At the same time, the tested ROP payloads are far from the complexity of a fully blown ROP-based implementation of the operations required for carrying out an attack, such as dumping a malicious executable on disk and executing it. Currently, this functionality is handled by the embedded shellcode, which in essence allows us to view these ROP payloads as sophisticated versions of return-to-libc. More complex ROP code will probably require a larger number of unique gadgets, some of them containing instructions that are currently not necessary, e.g., for directly invoking system calls. Given that even a single broken gadget is enough to render ROP code ineffective, this would increase the potential of in-place code randomization.

In any case, in-place code randomization raises the bar for the attacker, and significantly complicates the construction of robust ROP code. We should stress that the randomization coverage of our prototype implementation is a lower bound for what would be possible using a more sophisticated code extraction method [51, 65]. In our future work, we also plan to relax some of the conservative assumptions that we have made in instruction reordering and register reassignment, using data flow analysis based on constant propagation.

Combining In-Place Code Randomization with Existing Techniques. Given its practically zero overhead and direct applicability on third-party executables, in-place code randomization can be readily combined with existing techniques to improve diversity and reduce overheads. For instance, compiler-level techniques against ROP attacks [47, 54] increase significantly the size of the generated code, and also affect the runtime overhead. Incorporating code randomization for eliminating some of the gadgets could offer savings in code expansion and runtime overheads. Our technique is also applicable in conjunction with randomization methods based on code block reordering [33, 16, 43], to further increase randomization entropy.

In contrast to the above techniques, which modify the structure of the code image of a program by rearranging blocks of code, instruction set randomization (ISR) [42, 13] alters the instruction set that is “understood” by the underlying system. Legitimate programs are translated to a randomly chosen instruction set, and run normally on top of a randomized execution environment that supports the chosen instruction set. Any foreign code injected within a running process as a result of an attack would fail to execute correctly, because the actual instruction set used is unknown to any external observer.

Although ISR can be applied for protecting against any form of code injection, it is not effective against attacks such as return-to-libc and return-oriented programming, which are based on the reuse of code that already exists in the address space of a vulnerable process—irrespective of the underlying instruction set. Conversely, in-place code randomization does not offer any protection against code injection attacks. Instruction set randomization breaks any assumptions about the instruction

set used by a running process, while in-place code randomization breaks any assumptions about the location (and potentially the outcome, in case of non-intended code fragments) of certain instruction sequences that already exist in the address space of a process. The two techniques are thus complementary, and can be used in tandem to protect against both code injection and ROP attacks.

In-place code randomization at the binary level is not applicable for software that performs self-checksumming or other runtime code integrity checks. Although not encountered in the tested applications, some third-party programs may use such checks for hindering reverse engineering. Similarly, packed executables cannot be modified directly. However, in most third-party applications, only the setup executable used for software distribution is packed, and after installation all extracted PE files are available for randomization.

9.9 Conclusion

The increasing number of exploits against Windows applications that rely on return-oriented programming to bypass exploit mitigations such as DEP and ASLR, necessitates the deployment of additional protection mechanisms that can harden imminently vulnerable third-party applications against these threats. Towards this goal, we have presented in-place code randomization, a technique that offers probabilistic protection against ROP attacks, by randomizing the code of third-party applications using various narrow-scope code transformations.

Our approach is practical: it can be applied directly on third-party executables without relying on debugging information, and does not introduce any runtime overhead. At the same time, it is effective: our experimental evaluation using in-the-wild ROP exploits and two automated ROP code construction toolkits shows that in-place code randomization can thwart ROP attacks against widely used applications, including Adobe Reader on Windows 7, and can prevent the automated generation of ROP code resistant to randomization. Our prototype implementation is publicly available, and as part of our future work, we plan to improve its randomization coverage using more advanced data flow analysis methods, and extend it to support ELF and 64-bit executables.

Availability

Our prototype implementation is publicly available at <http://nsl.cs.columbia.edu/projects/orp>

Acknowledgements

We are grateful to the authors of Q for making it available to us, and especially to Edward Schwartz for his assistance. We also thank Úlfar Erlingsson and Periklis Akritidis for their valuable feedback. This work was supported by DARPA and the US Air Force through Contracts DARPA-FA8750-10-2-0253 and AFRL-FA8650-10-C-7024, respectively, and by the FP7-PEOPLE-2009-IOF project MALCODE, funded by the European Commission under Grant Agreement No. 254116. Any opinions, findings, conclusions, or recommendations expressed herein are those of the authors, and do not necessarily reflect those of the US Government, DARPA, or the Air Force.

References

1. Adobe CoolType SING Table “uniqueName” Stack Buffer Overflow. <http://www.exploit-db.com/exploits/16619/>.
2. Immunity Debugger. <http://www.immunityinc.com/products-immdbg.shtml>.
3. Integard Pro 2.2.0.9026 (Win7 ROP-Code Metasploit Module). <http://www.exploit-db.com/exploits/15016/>.
4. MPlayer (r33064 Lite) Buffer Overflow + ROP exploit. <http://www.exploit-db.com/exploits/17124/>.
5. /ORDER (put functions in order). <http://msdn.microsoft.com/en-us/library/00kh39zz.aspx>.
6. Profile-guided optimizations. <http://msdn.microsoft.com/en-us/library/e7k32f4k.aspx>.
7. Syzygy - profile guided, post-link executable reordering. <http://code.google.com/p/sawbuck/wiki/SyzygyDesign>.
8. White Phosphorus Exploit Pack. <http://www.whitephosphorus.org/>.
9. Wine. <http://www.winehq.org>.
10. *Intel 64 and IA-32 Architectures Software Developer’s Manual*. Volume 2 (2A & 2B): Instruction Set Reference, A-Z. 2011. <http://www.intel.com/Assets/PDF/manual/325383.pdf>.
11. M. Abadi, M. Budiu, U. Erlingsson, and J. Ligatti. Control-flow integrity. In *Proceedings of the 12th ACM conference on Computer and Communications Security (CCS)*, 2005.
12. A. V. Aho, M. S. Lam, R. Sethi, and J. D. Ullman. *Compilers: Principles, Techniques, and Tools (2nd Edition)*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 2006.
13. E. G. Barrantes, D. H. Ackley, T. S. Palmer, D. Stefanovic, and D. D. Zovi. Randomized instruction set emulation to disrupt binary code injection attacks. In *Proceedings of the 10th ACM conference on Computer and Communications Security (CCS)*, 2003.
14. K. Baumgartner. The ROP pack. In *Proceedings of the 20th Virus Bulletin International Conference (VB)*, 2010.
15. E. Bhatkar, D. C. Duvarney, and R. Sekar. Address obfuscation: an efficient approach to combat a broad range of memory error exploits. In *In Proceedings of the 12th USENIX Security Symposium*, 2003.
16. S. Bhatkar, R. Sekar, and D. C. DuVarney. Efficient techniques for comprehensive protection from memory error exploits. In *Proceedings of the 14th USENIX Security Symposium*, August 2005.
17. T. Bletsch, X. Jiang, V. Freeh, and Z. Liang. Jump-oriented programming: A new class of code-reuse attack. In *Proceedings of the 6th Symposium on Information, Computer and Communications Security (ASIACCS)*, 2011.

18. F. Bouchez. *A Study of Spilling and Coalescing in Register Allocation as Two Separate Phases*. PhD thesis, École normale supérieure de Lyon, April 2009.
19. E. Buchanan, R. Roemer, H. Shacham, and S. Savage. When good instructions go bad: generalizing return-oriented programming to RISC. In *Proceedings of the 15th ACM conference on Computer and Communications Security (CCS)*, 2008.
20. S. Checkoway, L. Davi, A. Dmitrienko, A.-R. Sadeghi, H. Shacham, and M. Winandy. Return-oriented programming without returns. In *Proceedings of the 17th ACM conference on Computer and Communications Security (CCS)*, 2010.
21. S. Checkoway, A. J. Feldman, B. Kantor, J. A. Halderman, E. W. Felten, and H. Shacham. Can DREs provide long-lasting security? the case of return-oriented programming and the AVC advantage. In *Proceedings of the 2009 conference on Electronic Voting Technology/Workshop on Trustworthy Elections (EVT/WOTE)*, 2009.
22. P. Chen, H. Xiao, X. Shen, X. Yin, B. Mao, and L. Xie. DROP: Detecting return-oriented programming malicious code. In *Proceedings of the 5th International Conference on Information Systems Security (ICISS)*, 2009.
23. F. B. Cohen. Operating system protection through program evolution. *Computers and Security*, 12:565–584, Oct. 1993.
24. Corelan Team. Corelan ROPdb. <https://www.corelan.be/index.php/security/corelan-ropdb/>.
25. Corelan Team. Mona. <http://redmine.corelan.be/projects/mona>.
26. L. Davi, A.-R. Sadeghi, and M. Winandy. Dynamic integrity measurement and attestation: towards defense against return-oriented programming attacks. In *Proceedings of the 2009 ACM workshop on Scalable Trusted Computing (STC)*, 2009.
27. L. Davi, A.-R. Sadeghi, and M. Winandy. ROPdefender: A practical protection tool to protect against return-oriented programming. In *Proceedings of the 6th Symposium on Information, Computer and Communications Security (ASIACCS)*, 2011.
28. S. Designer. Getting around non-executable stack (and fix). <http://seclists.org/bugtraq/1997/Aug/63>.
29. T. Dullien, T. Kornau, and R.-P. Weinmann. A framework for automated architecture-independent gadget search. In *Proceedings of the 4th USENIX Workshop on Offensive Technologies (WOOT)*, 2010.
30. R. El-Khalil and A. D. Keromytis. Hydan: Hiding information in program binaries. In *Proceedings of the International Conference on Information and Communications Security (ICICS)*, 2004.
31. Ú. Erlingsson. Low-level software security: Attack and defenses. Technical Report MSR-TR-07-153, Microsoft Research, 2007. <http://research.microsoft.com/pubs/64363/tr-2007-153.pdf>.
32. A. Fog. Calling conventions for different C++ compilers and operating systems. http://agner.org/optimize/calling_conventions.pdf.
33. S. Forrest, A. Somayaji, and D. Ackley. Building diverse computer systems. In *Proceedings of the 6th Workshop on Hot Topics in Operating Systems (HotOS-VI)*, 1997.
34. G. Fresi Roglia, L. Martignoni, R. Paleari, and D. Bruschi. Surgically returning to randomized lib(c). In *Proceedings of the 25th Annual Computer Security Applications Conference (ACSAC)*, 2009.
35. I. Guilfanov. Jump tables. <http://www.hexblog.com/?p=68>.
36. I. Guilfanov. Decompilers and beyond. Black Hat USA, 2008.
37. L. C. Harris and B. P. Miller. Practical analysis of stripped binary code. *SIGARCH Comput. Archit. News*, 33:63–68, December 2005.
38. Hex-Rays. IDA Pro Disassembler. <http://www.hex-rays.com/idapro/>.
39. X. Hu, T.-c. Chiueh, and K. G. Shin. Large-scale malware indexing using function-call graphs. In *Proceedings of the 16th ACM conference on Computer and Communications Security (CCS)*, 2009.
40. R. Hund, T. Holz, and F. C. Freiling. Return-oriented rootkits: bypassing kernel code integrity protection mechanisms. In *Proceedings of the 18th USENIX Security Symposium*, 2009.

41. R. Johnson. A castle made of sand: Adobe Reader X sandbox. CanSecWest, 2011.
42. G. S. Kc, A. D. Keromytis, and V. Prevelakis. Countering code-injection attacks with instruction-set randomization. In *Proceedings of the 10th ACM conference on Computer and Communications Security (CCS)*, 2003.
43. C. Kil, J. Jun, C. Bookholt, J. Xu, and P. Ning. Address space layout permutation (ASLP): Towards fine-grained randomization of commodity software. In *Proceedings of the 22nd Annual Computer Security Applications Conference (ACSAC)*, 2006.
44. S. Krahmer. x86-64 buffer overflow exploits and the borrowed code chunks exploitation technique. <http://www.suse.de/~krahmer/no-nx.pdf>.
45. C. Kruegel, W. Robertson, F. Valeur, and G. Vigna. Static disassembly of obfuscated binaries. In *Proceedings of the 13th USENIX Security Symposium*, 2004.
46. H. Li. Understanding and exploiting Flash ActionScript vulnerabilities. CanSecWest, 2011.
47. J. Li, Z. Wang, X. Jiang, M. Grace, and S. Bahram. Defeating return-oriented rootkits with “return-less” kernels. In *Proceedings of the 5th European conference on Computer Systems (EuroSys)*, 2010.
48. Microsoft. Enhanced Mitigation Experience Toolkit v2.1. <http://www.microsoft.com/download/en/details.aspx?id=1677>.
49. M. Miller, T. Burrell, and M. Howard. Mitigating software vulnerabilities, July 2011. <http://www.microsoft.com/download/en/details.aspx?displaylang=en&id=26788>.
50. S. S. Muchnick. *Advanced compiler design and implementation*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1997.
51. S. Nanda, W. Li, L.-C. Lam, and T.-c. Chiueh. Bird: Binary interpretation using runtime disassembly. In *Proceedings of the International Symposium on Code Generation and Optimization (CGO)*, 2006.
52. Nergal. The advanced return-into-lib(c) exploits: PaX case study. *Phrack*, 11(58), Dec. 2001.
53. T. Newsham. Non-exec stack, 2000. <http://seclists.org/bugtraq/2000/May/90>.
54. K. Onarlioglu, L. Bilge, A. Lanzi, D. Balzarotti, and E. Kirda. G-Free: defeating return-oriented programming through gadget-less binaries. In *Proceedings of the 26th Annual Computer Security Applications Conference (ACSAC)*, 2010.
55. V. Pappas, M. Polychronakis, and A. D. Keromytis. Smashing the gadgets: Hindering return-oriented programming using in-place code randomization. In *Proceedings of the 33rd IEEE Symposium on Security & Privacy (S&P)*, May 2012.
56. M. Parkour. An overview of exploit packs (update 9) April 5 2011. <http://contagiodump.blogspot.com/2010/06/overview-of-exploit-packs-update.html>.
57. M. Pietrek. An in-depth look into the Win32 portable executable file format, part 2. <http://msdn.microsoft.com/en-us/magazine/cc301808.aspx>.
58. P. Saxena, R. Sekar, and V. Puranik. Efficient fine-grained binary instrumentation with applications to taint-tracking. In *Proceedings of the 6th annual IEEE/ACM international symposium on Code Generation and Optimization (CGO)*, 2008.
59. E. J. Schwartz, T. Avgerinos, and D. Brumley. Q: Exploit hardening made easy. In *Proceedings of the 20th USENIX Security Symposium*, 2011.
60. F. J. Serna. CVE-2012-0769: the case of the perfect info leak, Apr. 2012. http://zhodiac.hispahack.com/my-stuff/security/Flash_ASLR_bypass.pdf.
61. H. Shacham. The geometry of innocent flesh on the bone: return-into-libc without function calls (on the x86). In *Proceedings of the 14th ACM conference on Computer and Communications Security (CCS)*, 2007.
62. H. Shacham, M. Page, B. Pfaff, E.-J. Goh, N. Modadugu, and D. Boneh. On the effectiveness of address-space randomization. In *Proceedings of the 11th ACM conference on Computer and Communications Security (CCS)*, 2004.
63. Skape. Lcreate: An anagram for relocate. *Uninformed*, 6, 2007.
64. Skape and Skywing. Bypassing Windows hardware-enforced DEP. *Uninformed*, 2, Sept. 2005.

65. M. Smithson, K. Anand, A. Kotha, K. Elwazeer, N. Giles, and R. Barua. Binary rewriting without relocation information. Technical report, University of Maryland, 2010. <http://www.ece.umd.edu/~barua/without-relocation-technical-report10.pdf>.
66. P. Solé. Defeating DEP, the Immunity Debugger way. <http://www.immunitysec.com/downloads/DEPLIB.pdf>.
67. P. Solé. Hanging on a ROPE. http://www.immunitysec.com/downloads/DEPLIB20_ekoparty.pdf.
68. P. Ször. *The Art of Computer Virus Research and Defense*. Addison-Wesley Professional, February 2005.
69. Y. L. Varol and D. Rotem. An algorithm to generate all topological sorting arrangements. *Comput. J.*, 24(1):83–84, 1981.
70. P. Vreugdenhil. Pwn2Own 2010 Windows 7 Internet Explorer 8 exploit. <http://vreugdenhilresearch.nl/Pwn2Own12010-Windows7-InternetExplorer8.pdf>.
71. D. A. D. Zovi. Mac OS X return-oriented exploitation. RECON, 2010.
72. D. A. D. Zovi. Practical return-oriented programming. SOURCE Boston, 2010.