Survey of Control-Flow Integrity Techniques for Embedded and Real-Time Embedded Systems

TANMAYA MISHRA, Virginia Polytechnic Institute and State University, USA THIDAPAT CHANTEM, Virginia Polytechnic Institute and State University, USA RYAN GERDES, Virginia Polytechnic Institute and State University, USA

Computing systems, including real-time embedded systems, are becoming increasingly connected to allow for more advanced and safer operation. Such embedded systems are resource-constrained, such as lower processing capabilities, as compared to general purpose computing systems like desktops or servers. However, allowing external interfaces to such embedded systems increases their exposure to attackers. With an increase in attacks against embedded systems ranging from home appliances to industrial control systems operating critical equipment that have hard real-time requirements, it is imperative that defense mechanisms be created that explicitly consider such resource and real-time constraints constraints. Control-flow integrity (CFI) is a family of defense mechanisms that prevent attackers from modifying the flow of execution. We survey CFI techniques, ranging from the basic to state-of-the-art, that are built for embedded systems and real-time embedded systems and find that there is a dearth, especially for real-time embedded systems, of CFI mechanisms. We then present open challenges to the community to help drive research in this domain.

Additional Key Words and Phrases: Survey, Control-Flow Integrity, Real-Time systems, Embedded Systems

1 INTRODUCTION

Today, computing systems communicate through a complex web of interconnections. For instance, the modern smart-phone can simultaneously capture photographs and videos at quality rivaling that of movie cameras, upload gigabytes of information to the internet, turn on lamps and automatically control thermostats, stream high fidelity music to the nearest speaker, and can even unlock a car. We now live in the age of internet-of-things (IoT [8]) where the physical world around us can be manipulated by a push of a button.

The convenience afforded by such interconnections is, unfortunately, countered by the inconvenience of dealing with malicious parties who try to take control of these connected devices to inflict monetarily, and in some cases, bodily harm. A simple smart bulb from a reputed company was exploited to launch a distributed denial-of-service (DDOS) attack [58]. While a DDOS attack may have, at the most, an economic impact on the victim, an attacker could reprogram the lights to blink so as to induce an epileptic attack in some individuals. Unfortunately, such instances of malicious behavior are not confined to small home appliances. Stuxnet [34] is a computer worm built to infect supervisory control and data acquisition (SCADA) systems. Infections of this worm were first uncovered in 2010, by when it had already infected nuclear reactor control systems and caused significant damage to Iran's nuclear program. Malicious entities could, theoretically, cause the reactors to fail and cause catastrophic damage to both life and property. Interestingly, many Stuxnet systems were air-gapped, i.e, did not have a direct connection to external systems. Instead, the infection spread from physical drives inserted by human operators.

However, the scope of system attacks and their related defenses has an extremely large body of prior work. Over the years, a variety of system defense mechanisms have been proposed for a wide range of threat models and system configurations. These mechanisms can be either hardware assisted [37], entirely in software [83], implemented in the system pre-deployment, such as compiler-based protections [66], detect attacks during system runtime [35]. Discussing the entire body of work of such defenses is beyond the scope of this survey. Here, we discuss a specific set of defense

1

mechanisms called control-flow integrity (CFI) for embedded, and particularly, real-time embedded systems. Our major contributions are:

- (1) We explore a number of recently proposed mechanisms targeting embedded systems, specifically those that are *resource-constrained*, such as reduced processing capabilities over general-purpose processing environments systems such as those found in desktop or server-grade equipment. Such embedded systems usually feature low-end processing environments such as microcontrollers (and their related underlying processor architecture) and identify key techniques that could provide inspiration for more robust real-time systems CFI design.
- (2) We find that there are very few CFI mechanisms built specifically for real-time embedded systems. Our exploration of the work for embedded systems show that there is an avenue to extend CFI techniques from general embedded systems to create powerful CFI mechanisms that uniquely leverage real-time requirements.
- (3) We consolidate our findings and present challenges and suggestions for future research in Section 6.

We give definitions for *embedded systems* and *real-time embedded systems* later in this section that defines the scope of our survey. We will now provide an overview of the type of attacks that are countered by CFI and an overview of CFI itself.

1.1 Scope of Attacks and Defenses: Control-Flow Attacks and CFI

To aid the discussion of CFI, which is the main focus of this survey, it is necessary to first describe the type and scope of attacks for which they are built. This family of attacks are collectively called control-flow attacks. We shall now discuss these types of attacks.

1.1.1 Control-Flow Attacks: Control-flow attacks capture and modify the flow of execution of a program. These attacks attack control information, that is information presented to a program during runtime that determines the path that a program takes to continue execution. A simple example of such information is the return address of a function call. See Figure 1.a for an example of the stack frame of a function call on a generic ARM architecture-based microcontroller. Here, the value stored in the LR field of the stack frame is popped into the special LR or link register. ARM calling convention [31], which is implemented by all compilers that officially support this architecture, utilizes the LR to implement the return sequence of a function call. Return sequences are implemented by branching on the LR such as by using the BX LR instruction. Therefore, the contents of the LR effectively constitutes control information. Control-flow attacks aim to modify such information to redirect program execution for malicious purposes. The same figure showcases a sample attack where the attacker utilizes a buffer-overflow bug in the code that writes to a memory buffer in the stack frame, such that it uses the bug to overwrite the LR information thereby tainting the return address with a desired target address. Therefore, when the function returns, the tainted value is popped and becomes the target of the branch statement. Note that since control-flow attacks redirect program execution, they are also sometimes called code redirection attacks.

Two broad categories exist in control-flow attacks. While each category is a large research domain by itself, we briefly describe them here for context:

(1) Code injection attacks - The sample attack we discuss above is a simple example of a code injection attack. As discussed above, the attack can be broken into two stages that are a) injecting (writing) code into some form of executable memory, followed by b) a redirection to the beginning of the injected code, such as by using the LR register. Due to code injection that takes place in stage a), these attacks are termed code injection attacks.

Code injection attacks have a large body of work [36, 40, 67]. However, such attacks have lost favor over time with advancements in software and hardware architecture. Note that an implicit assumption of the attack is that code is injected into executable memory, that is, the stack is executable. Therefore, to defeat such attacks, it is sufficient to introduce countermeasures that ensure that writeable memory addresses are not executable. A large body of research has been presented to counter code injection attacks with relatively inexpensive performance overheads [48, 53]. Even for lower-end processors, such as microcontrollers from the ARM Cortex-M family, prior work have implemented defenses [52]. Modern hardware now include architectural features such as the memory protection unit (MPU) that make it trivial for system designers to implement writeable but non-executable memory, an important requirement for code injection attacks to propagate [22]. Since such attacks can be defended against relatively easily, such attacks are outside the scope of this survey.

(2) Code reuse attacks - With the addition of defenses against a code injection attacks, a new class of attacks emerged that are collectively called code reuse attacks. These attacks are a logical extension of code injection attacks where attackers modify control information to reuse arbitrary sequences of code already present within the program binary to perform malicious operations. One of the most famous examples of code reuse attacks is the the return-oriented programming or ROP [70, 75, 84] attack. We provide more details of ROP, and defenses against such attacks, in Section 3. Increasingly sophisticated variants of ROP [17], such as some that do not even require return sequences [18, 28] have been proposed over the past decade. A large set of defenses have also been proposed to counter ROP and related attacks [20, 49, 80], showcasing the relevance and danger such type of attacks represent for modern systems.

Since control-flow attacks modify the control-flow of the program, it is necessary to maintain the *integrity* of the control-flow by detecting malicious control-flow deviation when it occurs. Therefore, any defense mechanism that enforces this integrity is called control-flow integrity (CFI) [4]. In this survey, we discuss CFI that is specifically designed to defend against code reuse attacks. Note that we alternatively refer to CFI as CFI enforcement, CFI mechanism, or CFI technique throughout this paper.

1.1.2 Control-Flow Integrity (CFI):. CFI is the set of system security techniques built to prevent an attacker from forcing a software system to execute code in an unintended manner. CFI focuses on ensuring that system code does not deviate from known software control-paths during system runtime. CFI mechanisms are built to address powerful threat models where it is assumed that the attacker can bypass all other defenses to infiltrate the system and force system software to execute in an arbitrary manner. There is a wealth of research in recent years that develop CFI mechanisms for increasingly complex and powerful attack scenarios [13, 30]. CFI mechanisms are also available in many commercial and production-grade software. For example, the Clang compiler implements control-flow violation detection mechanisms [1], and Microsoft has its own CFI implementation called Control Flow Guard [57] for its Windows operating system which has been available since Windows 8.1.

While the literature concerning CFI mechanisms (and techniques to bypass them [16, 32]) is rich with studies regarding the non-negligible performance and/or memory overhead of the mechanisms, few are built specifically for embedded systems and even fewer explicitly consider the real-time requirements of such systems. Therefore, we shall first look at CFI mechanisms for general embedded systems and then move towards mechanisms built explicitly for those with real-time constraints. We shall look at both software-based and hardware-assisted mechanisms, as well as a mechanism that takes advantage of the predictability of real-time systems. However, before we begin discussion of CFI techniques, we shall now define resource and real-time constraints.

1.2 Systems Considered: Embedded and Real-Time Embedded Systems

There exists numerous prior work that are excellent surveys and compilations regarding CFI defenses for general systems [4, 13, 30, 73, 78]. However none of these work explicitly consider system capabilities and constraints.. We now define the types of systems that we consider for the rest of this work, and their constraints that influence the design of CFI for such systems.

1.2.1 Embedded Systems: As discussed earlier, the Stuxnet worm was built specifically to target and control SCADA systems. A SCADA system is usually composed of a number of embedded computing systems built for specific operations, such as data gathering and actuator control. However, embedded computing systems themselves can be found in a wide variety of operating environments, ranging from complex SCADA systems, to robots used for medical procedures as well as small household appliances. These embedded systems are usually severely resource-constrained to minimize size, weight and power (SWaP), cost and/or simplify operations. Typically, they consist of microcontrollers that are low-end processors with integrated memory, executing software built to perform specific operations in a deterministic and predictable manner. For example, the modern vehicle can have over a hundred individual computing units, called Electronic Control Units (ECU) that control different functionalities of the vehicle. These units usually consist of a microcontrollers [51] that operate at a clock frequency an order of magnitude lower than the processors found in modern internet servers, and have similarly small amounts memory for storage and operation. These computing units control vehicle operations ranging from non-critical infotainment systems, to extremely critical Advanced Driver Assistance Systems (ADAS), such as anti-lock braking systems, whose failure could result in passenger loss-of-life. Further, the software for such systems may not be regularly updated due to the inaccessibility of their deployment locations. Therefore, once security vulnerabilities in the software are found, they may not be easily patched, making them lucrative targets for malicious entities. In addition, in the case of modern vehicles, increasing inter-vehicular connectivity to improve ADAS as well as increasing number of interfaces such as WiFi and Bluetooth for passenger convenience, has widened the attack surface that can be exploited by such entities [19]. Therefore, due to the wide range of applications of resource-constrained embedded systems and their increasing attack surface due to system inter-connectivity, it is imperative that such systems have built-in defense mechanisms to prevent their exploitation by attackers.

To summarize, our definition of embedded system is in the broader sense. That is, our definition encompasses embedded systems with fixed system resources (memory, processor, peripherals, etc.) where processing elements are embedded off-the-shelf microcontroller architectures such as ARM Cortex-M and ARM Cortex-R [90] or bespoke architectures that evolve from those that could be utilized in similar systems. Such processing environments have fewer architectural features than desktop or server grade processors and usually paired with slower/limited memory and peripherals for managing costs and/or special memory systems for redundancy and safety. Such systems are usually deployed in mission-specific applications in a wide range of domains, such as industrial, automotive, space and medical systems or even internet-of-things (IoT) systems. Our definition of such systems is broad since it allows us some flexibility to look at CFI mechanisms that may work for a specific type of embedded system, but could be applied to similar architectures with some modifications, giving us a broader field-of-view of the domain. For each mechanism we take a closer look at in later sections, we state the specific architectural considerations that informed its design. Note that for completeness we also briefly discuss some techniques that use external processing resources such in Section 4.5 and show their fundamental similarities with techniques that do not require external processing resources. However, we do not present in-depth information for these techniques since they utilize external processing resources that makes

it difficult to compare with techniques that do not require such external resources. Note that our definition of embedded systems assumes that such systems are resource-constrained and we interchangeably refer to embedded systems as embedded systems or resource-constrained embedded systems throughout this work.

1.2.2 Real-time embedded systems: Similarly, many such resource-constrained embedded systems require real-time guarantees. In the case of ADAS systems such as anti-lock braking systems, for example, multiple control loops (including actuator control) must be completed per second to maintain safe vehicle operation. We term such embedded systems as real-time embedded systems. If such a system misses any deadline, regardless of the correctness of the computation, the consequences could include the loss of life. When such guarantees are required atop resource-constraints, developing defense mechanisms for such systems become especially challenging.

We therefore focus on defense mechanisms that are built for embedded systems and real-time embedded systems. Since such systems have both *resource* and *real-time* constraints, considering systems that have a combination of these two types of constraints leads to a unique set of problems for designing useful CFI mechanisms for such systems. In general, some of the problems are:

- (1) Weaker processing capabilities as compared to general-purpose desktop or server grade systems constrains the complexity of the design and scope of the CFI mechanism that can be introduced in the system. Complex CFI would introduce unmanageable overheads that would break the real-time guarantees of the system. For example most of the defenses we discuss specifically for embedded systems in Section 4 detect irregularities in branch source and targets, individually for each branch. However, general-purpose architectures have more complex mechanisms available [20, 25] since such systems are not constrained by real-time guarantees and can accept greater performance reductions for higher degree of security.
- (2) In addition, due to reduced hardware capabilities, certain defense mechanisms that are built for general-purpose systems may not be directly applicable to resource-constrained embedded systems. For example some defenses [20] require advanced memory management features such as virtual memory which are not available on low-end microcontrollers. Therefore defenses for such systems require hardware/software workarounds to maintain acceptable levels of defense without hampering real-time operation.
- (3) Real-time systems require a study of the increased overhead due to CFI mechanisms, and its impact on the schedulability of the system. Many CFI techniques designed for general-purpose and, in fact, as we see later in Section 4, resource-constrained embedded systems do not discuss schedulability, nor do they discuss possible security-schedulability trade-offs that may be required to balance timing and security.
- (4) Other system parameters, such as power consumption are rarely considered when discussing CFI. Many resource-constrained embedded systems may have access to limited (such as battery-based) or intermittent (such as via renewable energy like solar) power supply. Such constraints are rarely discussed by prior work. We, in fact, realize a gap in knowledge with respect to impact of CFI and power consumption and suggest readers to explore this domain in future work (see Section 6).

To the best of our knowledge, this survey is the first to identify a gap in research of CFI mechanisms for real-time embedded systems, and propose future research avenues that could be considered by the real-time systems community. In this survey, we discuss CFI for real-time embedded systems and not general real-time systems that do not consider resource-constraints (such as memory or low end computation environments) typical to embedded systems. This is due to a lack of CFI literature that explicitly considers real-time constraints without considering resource constraints. On the other hand, we believe that our discussion of CFI for real-time embedded systems provides adequate coverage of

5

possible techniques that can be utilized, without many modifications, for any general real-time system. We also believe there is ample opportunity to investigate the unique hardware-software constraints of resource-constrained embedded systems and utilize real-time execution characteristics to aid the development of CFI techniques which are equally applicable to real-time systems that do not suffer from resource-constraints. Such timing based co-design, as we show in later sections, is severely lacking and we present a few possible paths of investigation for the reader to follow for future work in Section 6.

2 PAPER ORGANIZATION

The rest of this work is divided into 4 major sections. These are:

- (1) CFI Techniques for Backward and Forward-Edges (Section 3) We discuss different CFI designs, from both theoretical and practical approach, for general-purpose systems. This section provides the reader a general overview of how state-of-the art CFI mechanisms, both basic and advanced, are usually designed and implemented.
- (2) CFI for Embedded Systems (Section 4) We discuss different types of CFI techniques built specifically for resource-constrained embedded systems. Please note our definition of embedded systems is provided in Section 1.2. As stated in Section 1.2, the nomenclature "embedded systems" and "resource-constrained embedded systems" are synonymous and interchangeably used depending on context for clarity.
- (3) CFI for Real-Time Embedded Systems (Section 5) We then discuss how real-time considerations play into the design of CFI for embedded systems. Four specific techniques are considered that explicitly consider real-time constraints and discuss schedulability-security trade-offs and/or schedulability analyses.
- (4) Summary and Open Challenges (Section 6) We summarize our discussion of different CFI techniques and discuss some challenges from a real-time perspective, and from overall resource-constrained embedded system perspective.

Table 1 provides a brief overview of the relevant sections where we discuss specific CFI techniques, especially for Section 3.2, Section 4 and Section 5.

3 CFI TECHNIQUES FOR BACKWARD AND FORWARD-EDGES

We shall now look at some general techniques that are used in many CFI mechanisms. We will first look at techniques developed to prevent an attacker from modifying return sequences of function calls (backward-edge) or modifying other points-of-interest, such as indirect branches/function calls (forward-edge). Techniques for the former are well established and extensively utilized in mechanisms for embedded systems and real-time embedded (Section 4 and Section 5). However, some recently proposed advanced techniques for forward-edge CFI have not yet been considered for real-time embedded systems and are highlighted in Table 1. Note that for this section and the rest of this paper, "performance overhead" and "overhead", unless stated otherwise, are synonymous and refer to the increase in the CPU cycles required due to the addition of the CFI mechanism into the system. "Memory overhead" refers to the increase in the total memory (code and data) required to implement the mechanism, unless otherwise specified. Unfortunately, not all prior work discussed in this survey utilized the same benchmarking software and hardware. Neither did they always report memory overheads. We present the information regarding overheads as it was presented in the original work. We only quantitatively compare different work if the overheads have been measured using the same combination

CFI	Forward-edge		Backward-	M - d 1 1 1 1 1	
Mechanisms	Fine-	Coarse-	edge	Mechanism highlights	
	grained	grained			
Advanced forward-edge techniques for general systems (Section 3.2):					
BBB-CFI [45]		✓	✓	Block-based enforcement - binary-only approach	
				without need for CFG	
PathArmor [79]	✓	✓	✓	Context-sensitivity - requires architectural support	
CFI for embedded systems (Section 4)					
Silhouette [89] (Section 4.1)		✓	✓	Uses shadow-stacks and labeling	
Control-flow (Section 4.2) locking [12]	1	1	✓	Lazy + shadow-stack replacement.	
locking [12]	V	V	V	Lazy + shadow-stack replacement.	
μ RAI [6],				Register-based CFI - shadow-stack replacement	
Zipper Stack [54], (Section 4.3)			✓	Interrupt-handling (µRAI)	
PACStack [55]				2 7 .	
CFI CaRE [62], TZmCFI [50] (Section 4.4)			/	ARM TrustZone based shadow-stack, nested	
TZmCFI [50]			V	interrupts Stronger threat models	
HCFI [21] (Section 4.4)			✓	New ISA that integrates shadow-stack	
				operations in processor pipeline	
CFI for real-time embedded systems (Section 5)					
RECFISH [82] (Section 5.1.1)	✓		√	Large-scale schedulability study of common	
				CFI techniques applied to an RTOS	
Improve schedulability (Section 5.1.2)			✓	Searching the number of task jobs that can have	
by reducing security [43] (Section 5.1.2)				CFI turned on to improve schedulability	
Timing-deviation [11] (Section 5.2.1)				Detects control-flow deviation by excess	
				computation time	
ECFI [5] (Section 5.2.2)	✓	✓	✓	CFI for hard-real time PLC code that detects	
				abnormal increase in execution	

Table 1. Table of contents of advanced forward-edge CFI techniques discussed in Section 3.2, CFI techniques for embedded systems discussed in Section 4, and CFI techniques for real-time embedded systems discussed in Section 5. Important highlights of each technique and degree of coarseness of forward-edge path deviations is discussed.

of hardware and software. That said, we try to provide a qualitative discussion when possible to aid the reader in determining the pros/cons of the CFI technique based on the values reported.

3.1 Backward-edge CFI techniques

The first step to any control-flow attack is infiltration. There must be some flaw in the system that can be exploited by an external attacker to begin a control-flow redirection. A very common software flaw is the buffer overflow. Due to the tight memory restrictions of embedded systems, and the flat memory model due to the lack of complex (and potentially expensive, from both economic and performance perspectives) memory management units, buffer overflow or stack overflow flaws are common in resource-constrained embedded systems since they are usually programmed using memory-unsafe languages such as C/C++ [76]. A simple example of such a flaw is a statically allocated array that is filled past its capacity. Imagine such a flaw exists within a function call of a driver code that handles user input from a keyboard. In the absence of proper memory management, such flaws can be easily exploited to overwrite adjacent locations within the function stack frame as seen in Figure 1(a). Of particular interest is the return address value in the stack frame. Overwriting the return address with a target address ensures that when the function returns, the code will

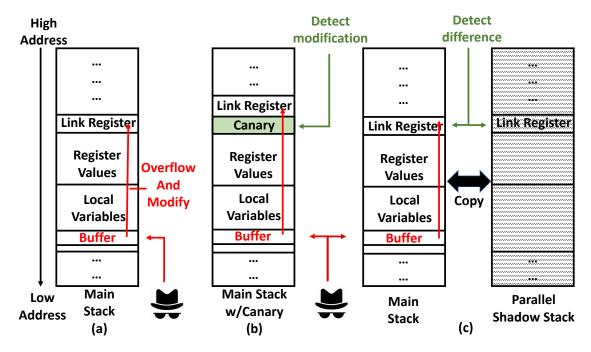


Fig. 1. (a) A function stack without any defenses, (b) Backward-edge CFI using an embedded canary, or (c) a parallel shadow stack (Section 3.1).

continue execution at the target address, successfully redirecting the flow of the program. The target address could be either a location within the pre-existing code memory or to some other memory address. A simple use case for the latter technique is to first *inject* the malicious code into the stack memory using the overflow vulnerability, and then set the return address to the start of the injected code. When the function returns, the injected code executes. Code-injection attacks can be thwarted with the help of memory protection mechanisms that implement the W \oplus X memory policy, i.e., prevent execution from writable memory. Such memory protections are now readily available in many commercial-off-the-shelf (COTS) low-end processors and microcontrollers. Therefore, the rest of our discussion will be focused on the consequences of the former technique of forcing the processor to continue execution at a target address in code memory.

Pointing the processor to an incorrect location by overwriting the return address is an example attack that serves as an entry point to a set of very powerful *code-reuse attacks*. For example, a well-studied sub-family of control-flow attacks is Return-Oriented Programming (ROP) [70]. A ROP attack is where an attacker chain together arbitrary code sequences (also called *gadgets*) that are already present on the device to achieve their objective. Post the seminal work by Shacham [75], ROP attacks have become increasingly popular and very sophisticated. It should be noted that using the return address to perform a control-flow diversion is also referred to as *backward-edge* control-flow attack. On the other hand, *forward-edge* control-flow attacks modify function pointers, or the targets of indirect function calls, to reuse code. An example is that by Checkoway et.al. [18] that modifies the target of indirect function calls to create gadget chains. Forward-edge defenses are discussed in the next section and are slightly more ambiguous in nature. It is interesting to note that all these attacks require exploiting an initial vulnerability such as a simple buffer overflow bug.

Two simple mechanisms to deal with backward-edge control-flow attacks are stack canaries [24] and shadow stacks [15]. Both these mechanisms, especially the latter, feature heavily in more sophisticated realistic CFI mechanisms for resource-constrained embedded systems. Stack canaries are special values inserted into the stack frame and are located in between the return address and the local statically allocated variables as seen in Figure 1(b). The concept behind using stack canaries is that an attacker overwriting the stack using a buffer overflow will have to first overwrite the canary value before overwriting the return address. Checking the canary value in the stack frame before a return operation can help determine whether the return address can be trusted. However, stack canaries can be bypassed by a sophisticated attacker, especially if the canary value is known (not random) or if the value can be guessed (not random enough). Further, they do not stop the attacker from overwriting local variables located before the canary value. By doing so, the attacker can still influence the function call operation [69].

Shadow stacks are a more sophisticated defense mechanism. Under the assumption that the attacker cannot access or modify a portion of the memory, a copy of the stack frames, or at least return addresses, is kept in that memory portion. Figure 1(c) presents an example of a shadow stack. This copy is updated during the initial stages of a function call (such as in the function prologue), and the return address is checked just before the return instruction is executed. If a discrepancy exists between the stored and actual addresses, it can be indicative of an attack. Shadow stacks are essentially more sophisticated canaries since both mechanisms indicate an attack by checking for discrepancies in the contents of the stack, with the major difference being that the shadow stack keeps a copy of the correct value [26]. While these mechanisms are relatively simple, applying them comes at a cost.

Dang et.al. [26] performed a study of the overheads caused by two different shadow stack implementations on the SPEC CPU2006 [2] standard suite of benchmarks on an x86 architecture processor. The first is a "traditional" shadow stack that has its own stack pointer and stores only the return addresses. The second is a "parallel" shadow stack that uses the same stack pointer as the main stack, however, the parallel shadow stack is stored at a different base address and records the return addresses while skipping over the other values in the stack frame (Figure 1(c)). Architecturally, this makes the parallel shadow stack faster than the traditional shadow stack since the same offset can be used for both the main and shadow stacks. The correct entry can be accessed by simply swapping out the contents of the stack base register which can be achieved with a single instruction. On the other hand, a traditional shadow stack would require additional code to maintain the stack as well as at least one extra instruction per operation to increment or decrement the shadow stack pointer for push and pop operations. Their measurements of the performance overhead shows that traditional shadow stack implementation, on average, introduces a 9.69% overhead (over a system without shadow stacks) while the parallel shadow stack introduces a 3.51% overhead. Worst case overheads of both were 52.5% and 19.6%, respectively. The cost of checking the return address was an additional 0.8%. On the other hand, stack canaries had an average performance overhead of 2.54%. At first glance the parallel shadow stack mechanism is clearly better suited to applications that are performance sensitive. As discussed above, the performance benefits of parallel shadow stacks is expected since accessing the relevant position in the parallel shadow stack only requires swapping the stack base register since both stacks share the same offset whereas multiple operations are required for an equivalent operation on traditional shadow stacks. However, the traditional shadow stack has its merits for a resource-constrained system with a low amount of memory.

3.2 Forward-edge CFI techniques

Forward-edge control-flow attacks are the logical extension to backward-edge attacks. The increasing popularity of backward-edge defense mechanisms forced attackers to consider other points-of-interest (POI) to redirect control-flow.

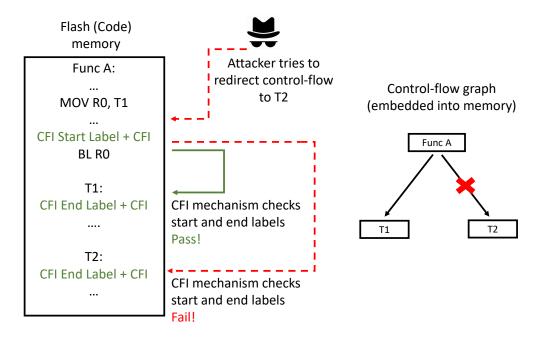


Fig. 2. Using labels and an embedded control-flow graph to enforce forward-edge CFI (Section 3.2).

These POI include indirect branches and indirect function calls via pointers. By attacking the destination of these branches, the attacker could call any arbitrary location without the need for return instructions [18].

Forward-edge CFI is difficult and, in general, subtler than backward-edge CFI. This is simply because *looking at the past is easier than predicting the future*. Forward-edge CFI techniques that could theoretically predict all possible combinations of branch start and end points are called *fine-grained* [4] CFI. Valid combinations of start and endpoints, essentially valid control-flow, can be represented as a *control-flow graph* (CFG). For example, Abadi et.al.'s [4] approach performs a binary static analysis using Vulcan [77] to generate a CFG and utilize said CFG to determine whether a branch is valid or not. A common mechanism to help enforce the valid control-flow paths in a CFG is *labeling*. Labeling is a process where all possible forward-edges that can be used by an attacker such as indirect branch locations, functions, and any other potential branch targets, are labeled with unique IDs. Figure 2 is an example of a labeling scheme where indirect branches and function prologues are labeled and matched against a CFG. When a branch occurs, the source label (such as an indirect branch) is checked against the destination label (such as a function) via code that has been instrumented into the binary (such as checks in a function prologue). A simple example of such an approach is presented in Figure 2.

An obvious problem of this approach, especially in the resource-constrained embedded systems, is the amount of memory required to store and enforce a CFG. However, more subtle issues arise in real-world cases. Many real-time embedded systems are industrial control systems, robotics systems, etc. In many cases, these environments run proprietary legacy software whose source code is difficult to obtain for analysis, or due to licensing issues, do not allow instrumentation. Due to these reasons, fine-grained CFG may not be possible to obtain, or the performance overhead associated with checking every branch may be prohibitive, especially in a real-time context. Therefore, many

coarse-grained CFI [88] have been proposed which allow varying degrees of relaxation of which branches or jumps need to be checked and which can be ignored. Due to reduced memory and processing requirements when utilizing coarse-grained CFG, coarse-grained forward-edge CFI are sometimes used for resource-constrained embedded systems. Due to the nature of coarse-grained CFI, such mechanisms may have blind spots that can be exploited by attackers [29]. A simple example is where a coarse-grained CFI allows any branch to any legal target, such as the start of a function, due to the unavailability of quality control-flow graphs. In such a case, the attacker could jump to targets which would have otherwise been identified as illegal by a fine-grained CFI. An interesting approach to overcome the need for a CFG, or the codebase to determine a CFG, is proposed by the authors of BBB-CFI [45] where the authors inspect the binary and divide it into *basic-blocks*, with each block having a single entry and exit point. A runtime mechanism prevents branches to the middle of a block, ensuring that the blocks are the smallest unit of code.

Interestingly, even fine-grained CFI can be defeated [16, 33], such as by exploiting the inability of current code static-analysis techniques to perfectly capture coding practices. Advanced forward-edge CFI techniques such as Van Der Veen et.al's PathArmor [79] can defend against such attacks. PathArmor logs control-flow transfers and then performs path verification by having access to the program CFG and performs a depth-first comparison of the logged transfers with the CFG to determine if the path taken during runtime is legitimate. This allows checking if a legitimate pair of source and destination addresses of a control-flow transfer are also contextually correct with respect to neighboring transfer events. However, the requirement for architectural support to record control-flow transfers prevent its direct application to low-end microcontroller-based systems that lack such specialized hardware.

3.3 A note on control-flow checking for soft errors

While CFI techniques are built considering an adversarial perspective, there exists a line of research that applies similar methodologies to detect erroneous control-flow redirection due to non-malicious soft-errors [41, 68, 74, 87]. These works utilize very similar techniques, such as by creating signatures for each basic block (code blocks that are delineated by control-flow transfers but do not contain any transfers themselves) and comparing currently executing basic block against a pre-determined graph of valid signature chains [63]. While such techniques utilize similar underlying principles to those discussed in prior sections, such as the forward-edge techniques in Section 3.2, soft-errors are generally one-shot errors that arise due to environmental factors. Control-flow redirection that is caused due to these errors are not easily predictable. For example, a redirection could take place due to reading the incorrect branch target from memory due to a bit-flip that took place in memory. However, control-flow redirection due to attacker control takes place under more predictable conditions (such as a buffer-overflow bug) and at a control-flow transfer point such as a branch/return statement. Further, advanced control-flow redirection, such as control-flow bending [16], where a control-flow transfer has valid start and end points but is incorrect only within the context of past control-flow transfers, cannot be detected by control-flow checking techniques since they are built to detect single-shot soft errors. For this survey, we will focus on techniques explicitly built for defending against various forms of control-flow redirection attacks.

4 CFI FOR EMBEDDED SYSTEMS

We now move towards more realistic CFI implementations in the context of resource-constrained embedded systems. The mechanisms presented here either combine techniques from Section 3 or propose entirely new techniques. Highlights of some of the mechanisms discussed in this section are presented in Table 1.

4.1 Implementation of basic techniques

We stated a pre-requisite in the prior section with respect to shadow stacks - ... Under the assumption that the attacker cannot access or modify a portion of the memory. This assumption does not have a straightforward justification in the context of embedded systems. As previously noted low-end embedded systems simply do not have complex memory management units to support well-known features such as virtual memory, which is now common in higher-end processors, let alone have special built-in mechanisms to support hiding shadow stacks from an attacker. Therefore, a successful CFI mechanism has to first wrangle the available hardware capabilities to support shadow stacks.

Zhou et.al's Silhouette [89] is an attempt to support shadow stacks on ARMv7-M [46], the architecture underlying ARM Cortex-M series of processors commonly found in embedded systems. It also supports forward-edge CFI checks. Silhouette, thus, is an example of how a sophisticated CFI mechanism would look like in the context of a resource-constrained embedded system. The ARMv7-M architecture supports two privilege levels in hardware, privileged and unprivileged. The optional memory protection unit (MPU) allows a system designer to decide access rights to an address. A limitation of the ARMv7-M architecture is that the MPU can be controlled by any privileged code. Most RTOS, such as FreeRTOS [10], by default, execute both the tasks and the operating system as privileged code to mitigate the overhead of switching privilege levels. This makes using the MPU to protect a shadow stack a moot point, simply because an attacker that has infiltrated the system, could re-program the MPU since they would most likely already execute under the privileged execution context.

Silhouette ensures that the MPU access rights are adhered to by working around this limitation. It replaces all store instructions, other than those that are supposed to directly store to the shadow stack, or the hardware abstraction layer (HAL) code, with unprivileged store variants, at compile time, to ensure adherence to the memory access policies defined in the MPU for the target address, regardless of the processor's current execution privilege level. The shadow stack is implemented in a similar manner as the parallel shadow stack explained in Section 3.1. To ensure that the store instructions with higher privilege levels are not abused by an attacker, Silhouette implements forward-edge CFI checks. Silhouette utilizes a labeling mechanism (Section 3.2) to guarantee forward-edge CFI [14].

On the performance front, Silhouette is benchmarked using well known embedded systems benchmark suites, namely CoreMark-Pro [23] and BEEBS [64]. We will see these same benchmarks being used in other approaches too in later sections, providing a common playing field. The maximum performance overhead reported for the two benchmark suites is 4.9% and 24.8%, respectively, and a code memory overhead of 8.9% and 2.3% respectively. The geometric mean of the performance overheads for all the benchmarks in each test suite is 1.3% and 3.4%, respectively. The approach used by Silhouette, which they term as *store hardening*, basically utilizes a memory management technique to hide the shadow stack from the attacker.

Another mechanism that can be used to prevent access to the shadow stack is called software fault isolation [59, 81] (SFI). SFI is a technique where the address space is partitioned into *fault domains*. Any code within a fault domain has unrestricted access to code or data within the same fault domain, but the partitioning scheme prevents the code from accessing any memory outside the fault domain. This is achieved by instrumenting load/store instructions during compile time to trigger the fault handler if the memory access takes place outside the fault domain. A variant of Silhouette is proposed that utilizes this technique by instrumenting store instructions to restrict them from writing to the shadow stack unless the store instruction is part of the shadow stack manipulation code. The authors note a higher performance overhead, with the geometric mean results being 2.2% and 10.2% respectively for the two benchmarks, which leads the authors to conclude that the store hardening approach is superior in performance. However, it would

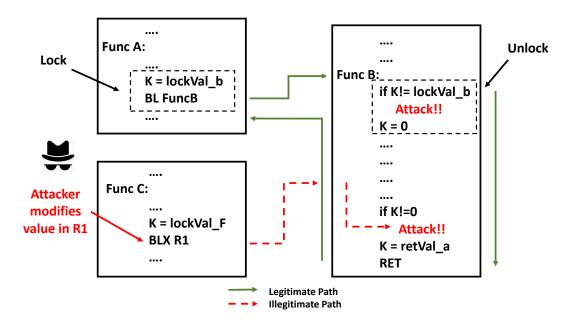


Fig. 3. Control-flow locking operation. Note the exclusive use of lock/unlocks for the entire operation (Section 4.2).

be interesting to note how the performance would vary if the shadow stack was protected using an approach similar to Aweke and Austin's [9] lightweight SFI for IoT systems that shows an overhead of just 1% on the MiBench [42] benchmarks. Their approach utilizes a small amount (150 lines) of trusted code that sets up the MPU to create the fault domains, trapping accesses outside the domain as memory access faults. Unfortunately, they do not present results using the CoreMark-Pro or BEEBS suites making direct comparisons difficult.

While the Silhouette and its variant provide a good overview of the well-known techniques of shadow-stacks and labels can be applied to a real low-end processor architecture, there are avenues to improve the operation of such systems.

4.2 Beyond the basics

While the techniques discussed in Section 3 consider forward-edge and backward-edge separately, some effort has been applied in recent years to develop more holistic mechanisms that apply to backward and forward-edges at the same time.

An example of such a mechanism is the Control-Flow Locking (CFL) technique [12]. This is also an example of a *lazy* CFI that trades-off attack detection speed with performance overhead. While CFL is not explicitly targeted at resource-constrained embedded systems, the mechanism can be implemented with similar memory and performance overhead as any general label-based CFI for detecting forward-edge control-flow attacks. CFL uses locks, instead of shadow stacks, to determine if an attacker has diverted control-flow to an arbitrary location. An overview of the CFL operation is given in Figure 3. The idea behind CFL approach is simple. Similar to how labels are generated based on the valid control-flow graph, *key values* are assigned to legitimate call/jump target locations. CFL targets indirect

calls/jumps as well as return instructions (an x86 architecture-based processor was assumed). Once the unique key values, which essentially represent valid edges in the control-flow graph, are generated, the authors propose to then instrument the target binary with instructions to lock and unlock control-flow paths using these key values. Every legitimate control-flow redirection start point, which may be an indirect call, jmp or ret instruction, is *preceded by a lock operation*, i.e., the key value is stored into a buffer. The assumption here is that the buffer is stored in a memory location such that it can be modified only by the lock and unlock subroutines, and not by attacker-controlled code. Once program control is redirected to a valid destination (such as a function entry-point), it is *immediately succeeded by an unlock operation* where the key value is validated, i.e., it is checked against a list of key values that could end up at this target location. If the values match, the key is zeroed out (*unlocked*) and execution continues as before. When the next control-flow redirection operation must take place, the key buffer is first checked to see if it contains a non-zero value. If it does, an attack is detected since no legitimate transfer would allow the key buffer to have a non-zero value due to the paired lock-unlock operations. Depending on the quality of the available CFG, this pairing of lock-unlock operations could be coarse or fine.

The overall mechanism is interesting due to its simplicity and the introduction of laziness. Not only does it prevent an illegitimate jump to a *valid* control-flow transfer site, but also automatically detects an illegitimate jump to an *invalid* control-flow site in recent history *without* requiring additional runtime memory such as using a shadow stack. Evaluations show that CFL can outperform fine-grained CFI mechanisms, with a maximum overhead of 21% v/s 31% overhead under Abadi et.al's [4] mechanism on the SPEC CPU2000 [2] benchmarks. The tradeoff? - As mentioned earlier the mechanism is lazy. The attacker can redirect control and can remain undetected until it is caught by the next locking site. While laziness allows the mechanism to work with the time and memory overhead similar to a labeling scheme, it could have interesting security repercussions especially in the context of the real-time embedded systems, many of which are used in industrial environments, controlling actuators in critical processes. If an attacker is able to send out control commands to these actuators before they are detected, the attacker can still inflict catastrophic damage. However, laziness is not inherently flawed. There is therefore an avenue to leverage real-time requirements to enforce timing bounds on laziness.

While CFL is an example of a CFI technique that re-purposes control-flow labels to solve both forward and backward control-flow attack detection at the same time, it still uses a form of memory protection. All the techniques discussed up to this point attempt to work around hardware limitations to enforce memory protection and are conservative. However, they do not take full advantage of the processor architecture or require radical software/hardware changes to improve performance.

4.3 Register-based shadow stacks

We will now discuss two approaches that would require significant software modifications to allow them to work. We will first briefly look at Zipper Stack [54] which is the more radical of the two since it proposes CPU architecture modifications to forego shadow stacks. The other is μ RAI [6] that is built for COTS embedded systems. It takes a more moderate approach by requiring reservation of parts of the CPU but can be implemented by recompiling the codebase with a modified compiler. Both implement backward-edge CFI.

Zipper stack aims to solve the problem of securing shadow stacks by replacing them with a set of processor architecture modifications. Shadow stacks, as discussed in Section 3 are inherently simple but require additional support to secure them from attacker manipulation. For example, Silhouette in Section 4.1 requires additional code instrumentation to secure the shadow stack Zipper Stack aims to solve this problem by replacing the shadow stack

with a single value stored in a special-purpose register called the *top register*. A separate register, the *key register*, holds a secret key. At the start of a new process, the key register and top register are initialized with random values. Each time a function call takes place, the top register is pushed onto the main stack alongside the actual return address. A message authentication code (MAC) algorithm, a cryptographic operation that is commonly used to authenticate messages from a known source, generates a new MAC from the top register value and the return address using the key in the key register. This newly created MAC is then stored in the top register. During a return sequence, the steps are reversed to authenticate the return address. First, the previous MAC value is popped from the stack and the MAC is recalculated using the return address and the popped MAC value. If the calculated MAC matches that currently in the top register, the return address is verified to be authentic. The processor replaces the top register with the popped value and continues execution at the return address. The purpose of the MAC based design is to reduce the attack surface. By utilizing the top register and chaining the MAC values with each successive function call, an attacker can only modify the return address and evade detection if it first modifies the value present in the top register (which is inaccessible to application code and is automatically updated by the hardware) before modifying the other MACs. Therefore, the rest of the MACs can be kept in non-secure memory that may be accessible to the attacker, reducing the amount of overhead introduced by accessing the "zipper stack" of MAC addresses.

The operation shows that Zipper Stack is heavily dependent on a) the efficacy of the MAC algorithm to ensure *collisions* (same MAC from different inputs) do not occur, b) the speed of the algorithm since every function call would constitute running the algorithm at least twice, and c) the attacker not being able to access the key register to forge MACs. For a), The authors use a well-known MAC algorithm, for b) the authors argue that a hardware implementation would allow MAC calculation in a single cycle, and for c) the authors argue that even if the key is leaked, the top register can only be modified at a call or a return operation. Their custom implementation on an FPGA with a RISC-V CPU achieves a 1.86% overhead on the SPEC CINT 2000 [2] benchmark.

While Zipper Stack presents a very radical approach that may never see wide-scale commercial adoption due to its hardware modifications, it is still interesting since custom architectures for specific applications, such as defense, are not uncommon in the embedded system world. In such cases, a custom architecture designed with optimized built-in defense mechanisms is not hard to envision. Interestingly, the use of MACs for authenticating return address may become possible very soon on commodity hardware. For example, PACStack [55] re-purposes the ARM pac instruction to create a MAC chain of return addresses, very similar to Zipper Stack. As part of the ARMv8.3-A PA extension, and soon to be available on SoCs based on ARMv8.3-A and later architecture revisions, pac allows generating pointer authentication codes (PAC) which are MACs generated on pointer values and stored along side the pointer. Similar to Zipper Stacks, the authors use a *chain register* to store PAC values which are generated from previous chain register values and the return address of a function call. When a return sequence takes place, similar to Zipper Stack, the reverse operation takes place. PACStack showed a geometric mean of 2.75% and 3.28% performance overhead on the SPECrate and SPECspeed (part of the SPEC CPU 2017 benchmark suite), respectively. PACStack provides a strong argument for MAC based shadow stack replacement, especially since it depends on architecture extensions which will soon be available in commodity hardware.

On the other hand, the authors of μ RAI take a similar but more realistic approach, especially on current-generation hardware. μ RAI is also concerned solely with the backward-edge, but instead of verifying the return address as is common with shadow stack approaches, μ RAI enforces Return Address Integrity (RAI) where the return address simply cannot be modified by an attacker. Their approach, in essence, is to prevent write access to the return address. μ RAI has the same set of requirements as many of the schemes we have discussed in previous sections, such as data execution

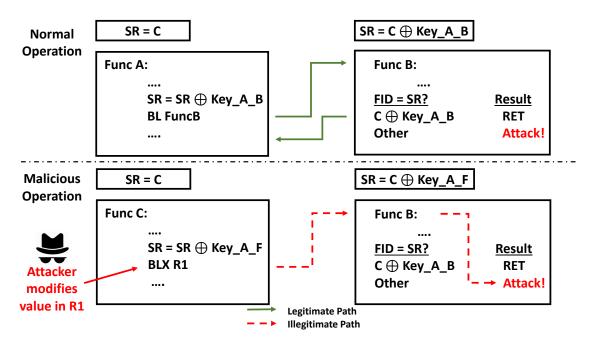


Fig. 4. µRAI operation. Shadow stack operation is implemented via SR register and FID table during return (Section 4.3).

prevention (DEP or $W \oplus X$) and an MPU. Similar to Zipper Stack, it requires that one of the processor registers is wholly dedicated to its operation and should never spill. This is called the *State Register* (SR). μ RAI's operation requires that the attacker cannot modify the register.

 μ RAI works by instrumenting code before branches and at return points, similar to CFL. It works solely with direct branches, i.e. branches with encoded destinations, and converts all indirect branches into direct branches by matching all possible start and endpoints. Figure 4 provides a basic overview of how μ RAI instrumented code looks like and operates. Every function *call site* is assigned a unique function key (FK). As is seen in the figure, a Function A can have multiple call sites to another Function B. μ RAI instruments code such that before every such call site, the value in the SR register is XOR'ed with the FK for the call site. This value is also called the Function ID (FID). The call goes through and Function B operates. At the point where Function B returns, it checks what the authors call the Function Lookup Table (FLT). This table has all the FIDs that could call this function. Based on which FID matches the value in the SR, the function returns to the corresponding location. Finally, the SR is XOR'ed with the same FK used before the branch, returning it to the original value before the function call. The authors tested their approach on an ARM Cortex-M4 based board and report a maximum performance overhead of 8.1% on the CoreMark [39] (a lighter variant of CoreMark-Pro) benchmark with an average of just 0.1%, making it comparable with shadow stack mechanisms discussed previously. However, it requires on average 34.6% extra flash memory for instrumentation and FLT.

The reader may have noticed that the possible return addresses are encoded into the code memory under DEP restrictions that prevent an attacker from modifying the code memory. DEP is enforced using the MPU. μ RAI, therefore, foregoes the return address that the processor may record in its stack, which is inherently writable memory, during a function call. Instead, it implements a function return mechanism that is implemented completely in code memory. This

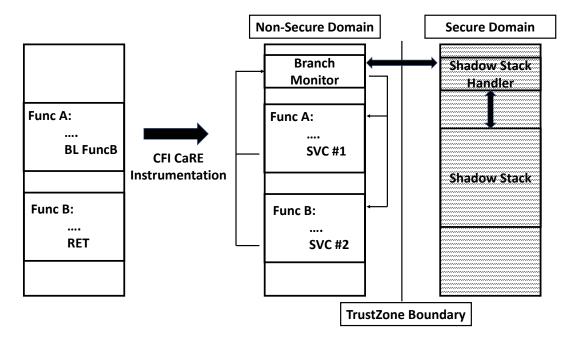


Fig. 5. CFI CaRE instrumented code. Branch and return instructions translate to SVCs (Section 4.4).

enforces μ RAI's goal of return address integrity. μ RAI is also the first mechanism that we have discussed in this survey, that explicitly considers interrupts. Since interrupts can occur at any time and can potentially interfere with shadow stack operations, they require explicit consideration. μ RAI instruments interrupt handler code to first save the return address which has been automatically stored on the stack by the hardware before the handler code is executed. μ RAI saves the return address to a safe memory hidden behind the MPU. Here μ RAI has to essentially create a shadow stack due to the limitation of the hardware. Supporting interrupts is a significant step to eventually supporting multi-threaded scheduling under a real-time operating system (RTOS). However, dedicating a register to μ RAI operations would require modifications to the compiler as well as incompatibility with embedded systems having a severely limited processing capacity, especially when the software requires large number of registers for computational purposes.

Unfortunately, none of these techniques improve forward-edge CFI. For example, in the case of μ RAI, the attacker could keep redirecting code execution using branch operations without allowing code to execute till an FID table. Therefore, such CFI mechanisms are helpful from only a performance or memory perspective over a regular shadow stack. That is, they do not provide any additional security guarantees, while requiring significant codebase changes or at least a modified compiler to support their operation.

4.4 CFI using processor architecture extensions

Before we finally move towards real-time aware CFI mechanisms, we will look at two mechanisms that depend on very modern processor architecture extensions such as ARM TrustZone [65]. TrustZone allows a processor to support two execution domains, *secure* and *non-secure*, each with its own address space with the secure domain having supervisory access to the non-secure domain. CFI designers have found creative ways to use it as part of their designs.

The first is Nyman et.al.'s CFI CaRE [62] that presents an alternative approach to secure the shadow stack, to that of Silhouette 4.1. An overview of its operation is given in Figure 5. While Silhouette uses binary instrumentation to prevent a privileged attacker from modifying the MPU that hides the shadow stack, CFI CaRE hides the shadow stack behind the TrustZone in the secure domain. CFI CaRE assumes that the original binary is only allowed to execute under the non-secure domain. It replaces all function calls with a *supervisory call* (SVC) that launches a special function called the *branch monitor*. The branch monitor runs in a privileged context and based on the parameter passed to the SVC that launches it, the branch monitor is able to identify if the source of the SVC is a branch or a return. It then calls secure domain code, passing the source identifier as a parameter, that updates the shadow stack. While the SVC ensures that all branches and returns are effectively trapped into the branch monitor, the TrustZone boundary ensures that non-secure domain code cannot view or modify the shadow stack. The authors used the Dhrystone (precursor to CoreMark) benchmarks to evaluate their work on an ARM Cortex-M23 processor. Performance overhead ranged between 13% to 513% with an overall 14.5% increase in flash memory consumption.

While CFI CaRE may seem like just a different implementation from previous approaches, it proposes a mechanism to address a crucial flaw in previous approaches with respect to embedded systems. The previously discussed approaches instrument binaries with no regard to the original layout. While this may be a non-issue for systems whose source code is available, many real-time embedded systems use proprietary legacy software and access to the source code may be limited. Further, due to memory and processor restrictions, these binaries are painstakingly built with strict adherence to page limits, available flash memory, etc. Unchecked binary instrumentation may destroy compatibility with the hardware. CFI CaRE's usage of SVC simply overwrites the branch or return instructions, keeping the original binary layout intact. However, it does require extra space for the branch monitor.

CFI CaRE also support interrupts and uses *trampolines* which are short sequences of code at the start of interrupt that call the secure domain to store the return address in a shadow stack. However, it does not support nested interrupts. If an attacker-controlled higher priority interrupt fires before the trampoline can store the return address in the shadow stack, the attacker-controlled interrupt code could rewrite the return address. When the lower priority interrupt finally gets to run, its trampoline would store a modified return address. Furthermore, nested interrupts can occur on an RTOS controlled system. For example, the timer tick could fire alongside interrupts from other peripherals. Kawada et.al.'s [50] TZmCFI fills in this gap. They too propose using the TrustZone to hide the shadow stack. However, they also extend the shadow stack concept to what they term as *exception shadow stacks* that support nested interrupts. They modify the trampolines such that every trampoline will complete all pending shadow stack transactions of lower priority interrupts before the interrupt body is allowed to execute. This ensures that if an attacker controls the interrupt body, it cannot affect the shadow stack copy of the interrupt return address. TZmCFI showed a performance overhead of up to 84% when supporting FreeRTOS as compared to FreeRTOS without CFI. For nested interrupts, the instrumented interrupts (with the trampolines) increased interrupt execution time from 30 cycles (un-instrumented) to 132-236 cycles, i.e., up to a 550% increase in execution time.

Other work that involves extending the architecture of the processing environment include work such as HCFI [21] suggest creating a new CFI enabled instruction set architecture (ISA) by modifying an existing ISA such as SparcV8's Leon3 [38]. They do so by adding new stages in the CPU pipeline to perform CFI operations such as shadow stack operations and show that performance overhead with respect to an umodified Leon3 core is less than 1% on their FPGA implementation for the SpecInt2000 benchmarks. While optimum performance can be achieved by designing a custom processor core as suggested here, unlike the TrustZone-based approaches discussed earlier, this would require significant investment to implement in real systems in the near future.

4.5 CFI using Separate Processing Environments

We wrap up our discussion of different CFI mechanisms for embedded systems with a brief note about CFI by utilizing off-chip processing environments since they behave very similarly to CFI achieved via TrustZone and utilize the same set of techniques presented in detail in Section 3. For example, techniques such as Abad et al.'s [3] uses a separate monitoring module to track the program counter and detect deviation from the control-flow. Similarly, SecMonQ [61] is designed for automotive systems and utilizes the Hardware Security Module (HSM) found in many commercial automotive ECU's to detect anomalous path behaviour. In a more general sense, techniques such as RTTV [85] utilize the Trusted Platform Module (TPM), a common co-processing environment used as a store for keys for cryptographic keys and perform a limited and static set of cryptographic operations in many embedded systems, can be used to store the CFG and perform regular measurements against the stored CFG. All these techniques inherit and apply the basic techniques presented in Section 3.

4.6 Section Summary

The techniques discussed in this section generally follow the basic techniques listed in Section 3. The proposed mechanisms either directly apply those basic techniques, or have progressively complex hardware modifications, from special registers to reduce the cost of shadow stacks (Section 4.3) to novel ISA (Section 4.4). However, the techniques do not inherently change the underlying principles of CFI and can be *conventional* by their nature. That is, they all verify the source and target destination addresses without much variation. Another important observation is that each of the techniques presented is uniquely tied to the underlying hardware for both performance and enforcement of CFI, making it difficult to compare their individual overheads. However, on a qualitative note, it is clear that the most performant CFI require radical hardware changes, such as integrating shadow stack operations into the pipeline of the processor [21].

A common theme in the techniques discussed, however, is the lack of any discussion regarding the implications of the overhead they introduce on systems where timing is critical, e.g. real-time systems. Real-time systems have certain characteristics that could be utilized to aid CFI and/or reduce the impact of the overhead introduced. We will now discuss these characteristics:

- (1) In periodic real-time systems, work is performed in a temporally predictable manner. That is, tasks execute during defined periodic intervals. CFI could utilize this predictable periodic nature to determine if an application is misbehaving due to attacker control.
- (2) The system is usually underutilized due to safety requirements. Since real-time systems are, in many cases, deployed in critical environments such as medical, industrial or automotive systems, such systems are designed to not perform work all the time to reduce or eliminate the possibility of missing deadlines. For example, the system is usually provisioned with enough computing resources such that tasks do not need to consume 100% of the computing resource at all times to complete by their deadlines. Therefore the system may have large periods of *slack* where the system idles, interspersed with heavy computation phases. CFI could utilize the slack thereby reducing localized spikes in computational load and reducing the possibility of missing deadlines. Note that although these systems may be underutilized, they are still considered to be resource-constrained. The underutilization is intentional due to safety concerns and any addition in the computational requirements must be done judiciously.

(3) The total system utilization at any given point of time is usually well characterized and there exist schedulability tests to determine if the system may be successfully scheduled without missing deadlines under a given scheduling algorithm and utilization. These tests may differ for different type of real-time task models (periodic tasks, aperiodic tasks, etc.). None of the techniques discuss their applicability and/or changes that must be introduced to satisfy these schedulability tests.

However, none of the techniques discussed in Section 4 consider timeliness, We now discuss CFI work that are specific to real-time embedded systems.

5 CFI FOR REAL-TIME EMBEDDED SYSTEMS

We have discussed multiple CFI techniques in the previous section for embedded systems. In this section, we survey the state-of-the-art mechanisms that consider real-time requirements. Unfortunately, there is little prior work that explicitly consider real-time properties of the system's operation. Therefore, this section discusses a few available CFI mechanisms. We divide our discussion into two parts, the first part covers techniques that are built specifically with an RTOS scheduler in mind, and the second discusses non-conventional CFI approaches. Highlights of the mechanisms discussed in this section are presented in Table 1.

5.1 CFI with an RTOS

5.1.1 An analytical approach for common CFI techniques. TZmCFI, presented in the previous section, is an example of CFI mechanisms for embedded systems that can work alongside an RTOS, or more specifically, a scheduler. A scheduler consists of supervisory code that decides when code that does actual work, i.e. complete the goal of the system, is able to run. A scheduler is critical to ensure system timeliness. While TZmCFI supports an RTOS, it lacks a study of system schedulability under different workloads. The recent work by Walls et.al. [82] addresses this deficiency in research. Their approach, called *RECFISH*, is an RTOS-aware CFI scheme. Since RECFISH shares several similarities with techniques discussed in prior sections, we will briefly discuss the mechanism and take a closer look at the evaluation results.

RECFISH is designed for ARM Cortex-R [56] processors that are built specifically for critical real-time applications. Like the Cortex-M series, they forego memory management units, and have special caching mechanisms to maintain predictability, support a small address space, but do not support TrustZone. RECFISH, instead, utilizes the MPU, like μ RAI, to enforce DEP. It assumes that the task code executes in the unprivileged mode while the RTOS runs in privileged mode. This ensures that if an attacker infiltrates a task, it cannot override the MPU settings. RECFISH is designed to be used with FreeRTOS and modifies it to allow setting up a per-task shadow stack (which only privileged code, such as the RTOS, can modify since it is hidden by the MPU), and modifies the scheduler to update the shadow stack when switching between tasks. Finally, RECFISH also instruments the binary to add labels to function prologues, as well as enforce shadow stack operations before (and after, in the function epilogue) the function body can execute. The labeling mechanism is used for enforcing forward-edge schemes, while the shadow stack operations are enforced by calling privileged shadow stack handling code using SVC just like that seen in CFI CaRE.

While the operation of RECFISH may look similar to multiple ideas presented in previous sections, the authors are the first to present a study of their approach's effect on real-time workloads. They evaluate and note a 21% performance overhead for their approach on the CoreMark benchmarks. Microbenchmarks show that RECFISH increases scheduler context switching time from 120 CPU cycles to 159. Further, the label checking and shadow stack operations increase

function prologue and epilogue overheads from 19 cycles (without any CFI operation) to 275 cycles. The authors then perform a large-scale schedulability study on simulated workloads. They randomly generated synthetic task sets with varying utilization values, task periods, and number of indirect branches. Utilization values ranging from 0.1% to 90% were considered. The overhead of task context switch (39 cycles) was incorporated into the task's worst-case execution time (WCET). For incorporating the function prologue and epilogue overheads (label checking for forward-edge CFI), the authors considered a varying number of indirect branches per task that were either 0, or ranging from 1 every $10^3 - 10^5$ cycles to 1 every $10^6 - 10^7$ cycles. Multiplying the number of branches with the 256 cycle overhead for the task yielded the overhead for the label checking mechanism which was then incorporated into the task WCET. RECFISH performs well for task sets where the number of tasks is few and each task has a high utilization, and when indirect branches are infrequent. However, the results show that up to 30% of the system utilization can become unusable for task sets with more frequent indirect branches and function calls and more tasks. Overall, RECFISH could schedule 85% of the 6 million task sets generated from 5760 different parameter combinations. The results show that well known CFI mechanisms such as shadow-stack and labeling could be used with a wide range of multi-threaded real-time workloads.

5.1.2 Trade-off security for schedulability. While RECFISH provides a schedulability study of common CFI techniques, Hao et al. [43] provides a novel technique to improve the schedulability of a real-time system by trading-off security with system schedulability. They focus on defending against ROP attacks (Section 1). They do so by selectively switching on CFI checks for a subset of instances (also called jobs) for each task in the system by exhaustively searching for the maximum set of jobs that can have CFI checks without hindering the schedulability of the system. The authors provide a comparison of an approximated scheduling algorithm that is designed to be faster to execute during runtime with respect to the exhaustive search algorithm which is determined to be optimal. Experimental results show that their approximation approaches optimality at lower (≤ 0.6) utilizations. A schedulability study shows that there is a sharp drop-off in schedulability of task sets if the CFI checks are added to task sets with utilization greater than 0.8. This observation echoes the results of the study of RECFISH that as task sets become "heavier", that is, have a higher utilization, schedulability sharply drops down to zero.

5.2 CFI utilizing timing deviations

5.2.1 Utilizing WCET: While RECFISH implements well known CFI techniques, Bellec et.al.'s [11] proposal utilizes the predictability of real-time systems to detect control-flow violations. An overview of the approach is provided in Figure 6. Their approach is based on the simple idea that an attacker will cause a control-flow violation to perform some malicious action. This will undoubtedly cause an increase in execution time, over and above the execution time of the system's tasks. Since real-time systems have well-defined task timing parameters, it is within reason to expect that an attacker-controlled execution would show a marked increase in execution time. A monitoring mechanism could, theoretically, detect such an increase and expose an attacker. The authors are able to support such a mechanism by first splitting the code base, consisting of a single task, into regions. Regions are either non-overlapping or located entirely within another region. Since the WCET of the task is known, each region within the task code is assigned a WCET of its own, called the maximal inner duration (MID). The MID of a region does not include the MID of a sub-region. Therefore, the sum of MIDs of all regions covering a task's code, would equal the task's WCET. The authors define another metric called the maximal attack window (MAW). For a set of monitored regions, the MAW is the maximum MID of that set. Therefore, the goal is to find the best possible set of regions such that a) the entire task code is covered, and b) the MAW is minimized. The authors perform a search, bounded by the available memory to store region boundaries as well

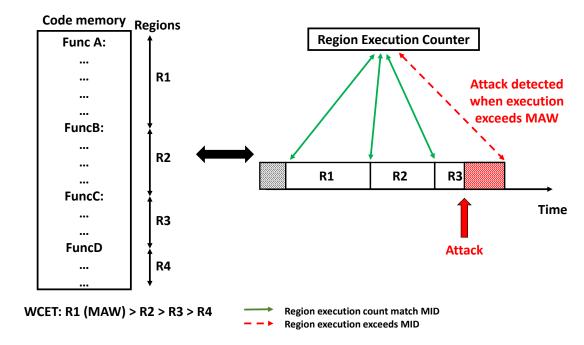


Fig. 6. Utilizing execution time (MID) as a metric to determine control-flow attacks (Section 5.2).

as performance metrics during runtime, to find the best possible set of regions. To evaluate their approach, the authors propose a custom hardware architecture that can detect when code execution enters and exits a region, as well as keep constant track of the time the processor spends within a region. If the time spent exceeds the MAW, an attack would be detected. The authors utilized two benchmark suites, Mälarden, and Polybench. They found that their approach had a mean latency of 95% (maximum of 99%) of the MAW before it detected an attack, where the MAW sizes ranged from few hundred up to over 160,000 CPU cycles. However, they found that their approach calculated MAWs of 600 or fewer CPU cycles for half of the benchmarks.

Due to the detection latency, this approach has similar issues as those that utilize laziness, specifically, the attacker could damage the system before it is detected. Further, it requires extensive modifications to the architecture to support it. However, it presents an interesting starting point for CFI mechanisms that effectively utilize the predictability of a real-time system to inform their approach.

5.2.2 Timing code in hard real-time context. We end our discussion of the state-of-the-art CFI for real-time systems with Abbasi et.al's [5] ECFI. ECFI is built for Programmable Logic Controllers (PLC) which are commonly found as the computing units for industrial-control systems. ECFI is a middle-ground approach, utilizing coarse-grained or fine-grained (depending on whether the code has pointer-based calls) CFI as well as exploiting the high predictability of the typical hard real-time system where PLCs serve as computational units, to detect if an attack causes a sudden increase in execution time to warrant the need to perform CFI checks. ECFI operates by capturing control-flow data in a global shadow-stack during system execution, and then check the data in a low-priority process. ECFI presents an amalgamation of traditional CFI techniques and utilization of predictability of the time-domain.

Note that there are related techniques to improve the schedulability of security mechanisms in general, such as Hasan et al.'s Contego framework [44] that introduces the concept of abstract security tasks into the system, but such techniques are not specifically designed for CFI and are not directly compatible with any of the work presented in this section.

5.3 Section Summary and Observations

Our discussion of CFI techniques for real-time embedded systems is summarized in Table 1. In general, we see a lack of techniques that consider timing constraints. While prior work has explored applying, with varying degrees, timing constraints to improving CFI schedulability there is still clear room for exploring this domain. For example, none of the techniques presented consider overloaded system conditions, or utilize timing to amortize the cost of CFI in such situations. For example, a periodic real-time system has well-defined intervals of slack. By deferring CFI operations to these slack intervals, it would be possible to reduce the effective in-line overhead that the CFI operation introduces while executing the system application, an observation we also state in Section 4.6. In our survey of CFI techniques for real-time systems, we have not found any technique that capitalizes on system slack in this manner. On the other hand, Hao et al.'s technique in Section 5.1.2 while useful to reduce the cost of CFI to maintain schedulability, can be considered incomplete in terms of security since only a subset of the code executed at runtime is actually checked. This could be exploited by a smart attacker, especially one aware of the technique used to decide which jobs do not have CFI checks. Some mitigation could be provided by randomizing the schedule using techniques such as using Yoon et al.'s TaskShuffler [86], but even such works have been shown to be defeated by carefully crafting an attack [60] that defeats the randomization. Essentially we do not see novel techniques that successfully use real-time constraints to amortize the cost of a complete implementation of CFI for real-time systems. Bellec et al's approach could be considered as a good starting point for creatively using timing constraints, however, it has its own failings which we discuss in Section 5.2.1.

6 SUMMARY AND OPEN CHALLENGES

For convenience, special terminology/mechanisms names that have been discussed before are listed here alongside the relevant section in the paper:

Silhouette - Section 4.1, Lazy - Section 4.2, Timing deviation - Section 5.2, BBB-CFI - Section 3.2, RECFISH, ECFI - Section 5.1.1, Context-sensitive - Section 3.2

A summary of our discussion in prior sections is presented in Table 2. Some common themes and omissions in the techniques presented are:

- (1) Most prior CFI work utilize some form of software-hardware bypass to accommodate hardware constraints present in resource-constrained embedded systems. The techniques trade-off performance and security to create the best possible compromise for their target hardware architecture.
- (2) The wide variety and heterogeneity of embedded system hardware make it difficult to all but qualitatively compare techniques in terms of memory and performance. Many require the use of custom/bespoke hardware architectures such as Zipper Stack (which requires a custom HMAC and special registers to speed up CFI). It is, therefore, difficult to judge if one technique is better. The applicability of any of the approaches we list for embedded and real-time embedded systems is dependent on the target application. We, therefore, only provide some qualitative discussion and summary, especially for the techniques discussed in-depth for embedded systems in Section 4 to aid the reader.

Category	Technique and Summary
Implementation: Standard CFI techniques on different architectures	Silhouette - Shadow stack and binary labeling on ARM Cortex-M RECFISH - Shadow stack and binary labeling on ARM Cortex-R
Design Changes: Non-standard CFI techniques utilizing standard control-flow start and end points	1) Control-Flow Locking - Lazy control-flow evaluation. Single technique for forward and backward-edge 2) uRAI - Collapse shadow stack into a single register using XOR operations. 3) Zipper Stack - Custom hardware to collapse shadow stack into a single register via HMAC operations.
Modern hardware architecture: Techniques that utilize new processor architecture features	1) CFI Care - Shadow stack hidden by ARM TrustZone 2) TZmCFI - Nested interrupts (RTOS) aware shadow stack in ARM TrustZone. 3) PACStack - ARM pointer authentication (ARMv8.3-A) utilized for collapsing shadow stack in single register
Underlying Principle: CFI techniques that detect control-flow deviations using non-standard principles	Timing deviation - Detect WCET violation of code segments using custom hardware ECFI - Built for PLCs. Detects timing violations code during runtime

Table 2. A summary of techniques discussed in depth in Section 4 and Section 5

- (3) With the exception of Bellec et al.'s work [11], the design of conventional CFI for embedded and/or real-time embedded systems can primarily be viewed as memory-based, where CFI is performed by detecting deviations from expected instruction memory accesses. There is no fundamental difference in the detection methodology across all the presented techniques.
- (4) Real-time CFI mechanisms, other than techniques such as that presented by Bellec et al.'s work [11], are *adhoc* in design. None of the techniques seem to utilize the strict timing requirements of the system to aid CFI. CFI, in essence, is detecting deviation in system behavior and timing critical systems depend heavily on being temporally correct. However, utilizing temporal guarantees exclusively to detect abnormal behavior can reduce the effectiveness of the mechanism as we discuss in Section 5.2 In fact, for a real-time embedded system, some assumptions can be made (we will discuss a possible approach later in this section) that can synergistically aid conventional CFI and improve its performance.

In this section we first present open challenges to the real-time community based on our understanding of the state-of-research in CFI for real-time embedded systems. We also present general consideration points that have not yet been incorporated into CFI designs.

6.1 Real-time challenges

We believe there exists two broad avenues of research that could be undertaken immediately, considering the state-ofthe-arts.

Bounded laziness: CFI designs for real-time systems are few in number and do not seem to capitalize on system predictability. In particular, laziness, such as that introduced by control-flow locking is a promising mechanism for hard real-time systems due to its ability to defer CFI checks. However, a drawback of their approach, and Bellec et.al's timing deviation based mechanism (Section 5.2) is the lack of expressiveness in the threat model, specifically, the time at which an attacker is able to affect the system. For example, the proposed mechanisms fail to consider that an attacker could modify and produce system outputs, such as sending messages via a network controller to other systems, before the CFI

mechanism detects an attack. On the other hand, conventional CFI techniques have an unnecessary sense of *urgency* since CFI is performed as close to when control-flow path changes as possible. For example, the mechanism presented in Silhouette adapts well-known CFI techniques which all perform CFI during a control-flow transfer event. We believe there is a middle-ground that can improve performance and still maintain the *usefulness* of CFI. That is, the purpose of CFI to detect an attacker before they are able to damage the system, is still maintained. This is because real-time systems inherently have discrete and well-known time instances where they must generate system output. For example, a typical control-system in an industrial environment, would command an actuator after a defined interval of time. In such instances, there is no need to *urgently* perform CFI, but CFI work can be deferred to much later by recording the control-flow transfer event and then verifying at a later stage before the actuator command is sent out. The advantage of such mechanisms would be a reduction in temporal overload situation which plagues current CFI implementation since they must be performed while system code is executing. However, the tradeoff is increased memory usage to record control-flow events. We believe there is ample opportunity to capture such memory-timing-security tradeoff situations in real-time systems due to the higher degree of predictability over general systems. Essentially, there is a need to define how lazy CFI can be, and develop system/task models that enforce these boundaries.

Multi-thread/core scheduling: RECFISH and ECFI showcases the applicability of well-known CFI techniques to multi-threaded hard real-time systems. We believe there is an opportunity to extend the concept of bounded laziness to multi-threaded systems and utilize the large pool of available real-time scheduling theory to tighten the bounds. For example, there could be an opportunity to steal system slack for performing CFI operations. In the case of multi-core scheduling, cores could be dedicated to performing CFI operations. From a security perspective, ECFI implicitly trusts the scheduler's integrity. However, in advanced threat models where an attacker could have the privilege to disrupt scheduler operations, such as modifying the system timer to warp the scheduler's sense of time, such defense mechanisms could fail. Prior work to secure time sources, such as TimeSeal [7] could provide some inspiration to solve this problem.

Determining CFI related workload attributes: Addressing the previous challenges would also require determining the real-time properties of CFI operations, such as the WCET of CFI operations, or how CFI operations would be incorporated into advanced real-time models such as those that consider varying task periods, varying number of real-time tasks, the effect of servicing aperiodic tasks, etc. The WCET of CFI too could be difficult to accurately determine especially if the mechanism operates on historical control-flow data, such as in context-sensitive CFI, where the amount of data can vary during system operation.

6.2 General challenges

In addition to the real-time system specific challenges listed above, there are some general considerations that should be incorporated into future designs. The following challenges are not just limited to CFI mechanisms but the system security research in general.

Power consumption: An often overlooked component of embedded system development is power consumption. This is also evident in every CFI design reviewed in this paper. None of the mechanisms consider power consumption, which is especially important in embedded systems operating off batteries and deployed in the field. Some designs such as that provided by Das et al. [27] provide power consumption measurements of their custom control-flow checking hardware design implemented on an FPGA. However, such measurements are an exception rather that the norm with respect to CFI research. Custom designs presented by other work such as Zipper Stack [54] do not provide information regarding power consumption, making it difficult to decide applicability of such work to severely power-constrained and hard

real-time environments such as heart pacemakers. Since CFI techniques such as shadow stacks have high memory access rates, impacts on system power consumption of different techniques must be considered.

We believe that, alongside real-time scheduling theory, techniques such as Dynamic Voltage Frequency Scaling (DVFS), backed by an extensive pool of scheduling algorithms that utilize DVFS [71, 72], could provide significant reduction in system power consumption and interesting schedulability issues. Interestingly, a logical correlation can be made between coarse-grained CFI and reduced power consumption, by virtue of reduction of CFI checks that are required due to the coarseness of the design. A study on the relation between coarse CFI and power consumption of design on commercial-off-the-shelf hardware could have an immediate impact within the research community, providing researchers guidance on which type of CFI and what aspects of CFI design have the worst effects on power consumption. We also believe that a new class of schedulability-power co-design problems could arise from utilizing laziness in CFI to limit the peak power consumption of a system by carefully differing CFI to low power consumption phases of the system.

The goal of CFI is similar to that of system reliability improvement techniques, i.e., to prevent incorrect execution and/or detect when incorrect execution occurs. There is a large amount of prior work that discuss mechanisms to implement recovery schemes with minimal impact to system power consumption. Such work could be used as inspiration to create energy-aware CFI mechanisms.

Portability: In general, CFI for resource-constrained embedded systems adapt well-known CFI techniques such as shadow stacks and labeling to such systems while working around their limitations. A primary observation is that many of these workarounds are very specific to the hardware platform that the authors target. For example, Silhouette targets ARMv7-M and therefore modifies store instructions to use the MPU on this architecture. Since these limitations are hardware-specific, designing realistic CFI mechanisms for such systems that are also portable is difficult. Unlike desktop or server-grade hardware, where commodity systems usually include processors with similar underlying architecture, embedded system utilize architectures from ARM, RISC, MIPS, etc. as well as application-specific designs. Designing a one-size-fits-all mechanism for such a wide-range of target architectures is a difficult challenge. Further, architectures such as ARM are very modular, allowing hardware vendors a high-degree of flexibility to add or remove features to adjust manufacturing costs and provide a wide portfolio of devices at every price point. There is, thus, a need to design feasible CFI mechanisms that operate completely in software (or with minimal hardware requirements), to allow for portable designs. However, the overhead of such designs remains to be seen.

Advanced CFI and beyond: As discussed in Section 3.2, there is a need to consider context-sensitivity in real-time embedded systems to thwart attacks that can bypass even fine-grained CFI. We are not aware of the existence of such techniques. Finally, there is a gap in research for embedded and real-time embedded systems regarding state-of-the-art data oriented programming [47] (DOP) attacks. These do not redirect control-flow, but attack program data, such as the counter variable used for a loop. Such attacks cannot be mitigated using any of the CFI designs discussed in this paper since they do not cause deviations in control-flow path. Note that techniques such as timing deviation detection discussed in Section 5.2 may be able to detect such attacks, but the assumption here is that the attacker is knowledgeable and does not violate the MAW during an attack. Data oriented attacks are powerful and have been shown to be capable of influencing program output as well as disclose private information.

7 CONCLUSION

We have examined multiple CFI schemes in the paper, starting from the core mechanisms that help enforce CFI, to the necessary workarounds required to support them in resource-constrained embedded environments. We have also looked

at the modifications necessary to support real-time schedulers and how real-time characteristics can be effectively utilized for CFI. While CFI has been adopted by higher-end systems, designs for resource-constrained embedded systems are still mostly academic and not yet widely deployed due to unmanageable performance overhead in some cases. As we have seen CFI will undoubtedly have overhead due to hardware constraints, but techniques such as laziness that trade-off detection speed with overhead could provide an interesting avenue for future work.

REFERENCES

- [1] 2020. Clang 12 documentation. Retrieved 2020-10-24 from https://clang.llvm.org/docs/ControlFlowIntegrity.html
- $[2]\ \ 2020.\ Standard\ Performance\ Evaluation\ Corporation.\ \ https://www.spec.org/benchmarks.html$
- [3] Fardin Abdi Taghi Abad, Joel Van Der Woude, Yi Lu, Stanley Bak, Marco Caccamo, Lui Sha, Renato Mancuso, and Sibin Mohan. 2013. On-chip control flow integrity check for real time embedded systems. In 2013 IEEE 1st International Conference on Cyber-Physical Systems, Networks, and Applications (CPSNA). IEEE, 26–31.
- [4] Martín Abadi, Mihai Budiu, Úlfar Erlingsson, and Jay Ligatti. 2009. Control-flow integrity principles, implementations, and applications. ACM Transactions on Information and System Security (TISSEC) 13, 1 (2009), 1–40.
- [5] Ali Abbasi, Thorsten Holz, Emmanuele Zambon, and Sandro Etalle. 2017. ECFI: Asynchronous Control Flow Integrity for Programmable Logic Controllers. In Proceedings of the 33rd Annual Computer Security Applications Conference.
- [6] N. S. Almakhdhub, Abraham A. Clements, S. Bagchi, and M. Payer. 2020. µRAI: Securing Embedded Systems with Return Address Integrity. In NDSS.
- [7] Fatima M Anwar, Luis Garcia, Xi Han, and Mani Srivastava. [n.d.]. Securing Time in Untrusted Operating Systems with TimeSeal. In 2019 IEEE Real-Time Systems Symposium (RTSS).
- [8] Luigi Atzori, Antonio Iera, and Giacomo Morabito. 2010. The internet of things: A survey. Computer networks 54, 15 (2010), 2787–2805.
- Z. B. Aweke and T. Austin. 2018. uSFI: Ultra-lightweight software fault isolation for IoT-class devices. In 2018 Design, Automation Test in Europe Conference Exhibition (DATE). 1015–1020. https://doi.org/10.23919/DATE.2018.8342161
- [10] Richard Barry et al. 2008. FreeRTOS. Internet, Oct (2008).
- [11] Nicolas Bellec, Simon Rokicki, and Isabelle Puaut. 2020. Attack detection through monitoring of timing deviations in embedded real-time systems. In ECRTS 2020-32nd Euromicro Conference on Real-Time Systems. 1–22.
- [12] Tyler Bletsch, Xuxian Jiang, and Vince Freeh. 2011. Mitigating Code-Reuse Attacks with Control-Flow Locking. In Proceedings of the 27th Annual Computer Security Applications Conference (Orlando, Florida, USA) (ACSAC '11). Association for Computing Machinery, New York, NY, USA, 353–362. https://doi.org/10.1145/2076732.2076783
- [13] Nathan Burow, Scott A Carr, Joseph Nash, Per Larsen, Michael Franz, Stefan Brunthaler, and Mathias Payer. 2017. Control-flow integrity: Precision, security, and performance. ACM Computing Surveys (CSUR) 50, 1 (2017), 1–33.
- [14] Nathan Burow, Scott A. Carr, Joseph Nash, Per Larsen, Michael Franz, Stefan Brunthaler, and Mathias Payer. 2017. Control-Flow Integrity: Precision, Security, and Performance. ACM Comput. Surv. 50, 1, Article 16 (April 2017), 33 pages. https://doi.org/10.1145/3054924
- [15] Nathan Burow, Xinping Zhang, and Mathias Payer. 2019. SoK: Shining light on shadow stacks. In 2019 IEEE Symposium on Security and Privacy (SP).
- [16] Nicholas Carlini, Antonio Barresi, Mathias Payer, David Wagner, and Thomas R. Gross. 2015. Control-Flow Bending: On the Effectiveness of Control-Flow Integrity. In 24th USENIX Security Symposium (USENIX Security 15). USENIX Association, Washington, D.C., 161–176. https://www.usenix.org/conference/usenixsecurity15/technical-sessions/presentation/carlini
- [17] Nicholas Carlini and David Wagner. 2014. {ROP} is still dangerous: Breaking modern defenses. In 23rd {USENIX} Security Symposium ({USENIX} Security 14). 385-399.
- [18] Stephen Checkoway, Lucas Davi, Alexandra Dmitrienko, Ahmad-Reza Sadeghi, Hovav Shacham, and Marcel Winandy. 2010. Return-oriented programming without returns. In Proceedings of the 17th ACM conference on Computer and communications security. 559–572.
- [19] Stephen Checkoway, Damon McCoy, Brian Kantor, Danny Anderson, Hovav Shacham, et al. 2011. Comprehensive experimental analyses of automotive attack surfaces.. In USENIX Security Symposium.
- [20] Yueqiang Cheng, Zongwei Zhou, Yu Miao, Xuhua Ding, and Robert H Deng. 2014. ROPecker: A generic and practical approach for defending against ROP attack. (2014).
- [21] Nick Christoulakis, George Christou, Elias Athanasopoulos, and Sotiris Ioannidis. 2016. HCFI: Hardware-enforced control-flow integrity. In Proceedings of the Sixth ACM Conference on Data and Application Security and Privacy. 38–49.
- [22] Abraham A Clements, Naif Saleh Almakhdhub, Khaled S Saab, Prashast Srivastava, Jinkyu Koo, Saurabh Bagchi, and Mathias Payer. 2017. Protecting bare-metal embedded systems with privilege overlays. In 2017 IEEE Symposium on Security and Privacy (SP). IEEE, 289–303.
- [23] EEMBC The Embedded Microprocessor Benchmark Consortium et al. 2015. CoreMark-Pro. Retrieved January 22 (2015), 2019.
- [24] Crispan Cowan, Calton Pu, Dave Maier, Jonathan Walpole, Peat Bakke, Steve Beattie, Aaron Grier, Perry Wagle, Qian Zhang, and Heather Hinton. 1998. Stackguard: Automatic adaptive detection and prevention of buffer-overflow attacks.. In USENIX security symposium, Vol. 98. San Antonio, TX, 63-78.

- [25] Stephen Crane, Christopher Liebchen, Andrei Homescu, Lucas Davi, Per Larsen, Ahmad-Reza Sadeghi, Stefan Brunthaler, and Michael Franz. 2015.
 Readactor: Practical Code Randomization Resilient to Memory Disclosure. 2015 IEEE Symposium on Security and Privacy (2015), 763–780.
- [26] Thurston HY Dang, Petros Maniatis, and David Wagner. 2015. The performance cost of shadow stacks and stack canaries. In Proceedings of the 10th ACM Symposium on Information. Computer and Communications Security. 555–566.
- [27] Sanjeev Das, Wei Zhang, and Yang Liu. 2016. A fine-grained control flow integrity approach against runtime memory attacks for embedded systems. IEEE Transactions on Very Large Scale Integration (VLSI) Systems 24, 11 (2016), 3193–3207.
- [28] Lucas Davi, Alexandra Dmitrienko, Ahmad-Reza Sadeghi, and Marcel Winandy. 2010. Return-oriented programming without returns on ARM. Technical Report. Technical Report HGI-TR-2010-002, Ruhr-University Bochum.
- [29] Lucas Davi, Ahmad-Reza Sadeghi, Daniel Lehmann, and Fabian Monrose. 2014. Stitching the Gadgets: On the Ineffectiveness of Coarse-Grained Control-Flow Integrity Protection. In 23rd USENIX Security Symposium (USENIX Security 14). USENIX Association, San Diego, CA, 401–416. https://www.usenix.org/conference/usenixsecurity14/technical-sessions/presentation/davi
- [30] Ruan de Clercq and Ingrid Verbauwhede. 2017. A survey of hardware-based control flow integrity (CFI). arXiv preprint arXiv:1706.07257 (2017).
- [31] R Earnshaw. 2005. ARM Procedure Call Standard for the ARM Architecture.
- [32] M. A. El-Sarraf, A. El-Sayed Abdo, and Mahmoud A. Abdulwahab. 2013. Usability of epoxy/ilmenite composite material as an attenuator for radiation and a restoration mortar for cracks. Annals of Nuclear Energy 60 (2013), 362–367.
- [33] Isaac Evans, Fan Long, Ulziibayar Otgonbaatar, Howard Shrobe, Martin Rinard, Hamed Okhravi, and Stelios Sidiroglou-Douskos. 2015. Control jujutsu: On the weaknesses of fine-grained control flow integrity. In Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security. 901–913.
- [34] Nicolas Falliere, Liam O Murchu, and Eric Chien. 2011. W32. stuxnet dossier. White paper, Symantec Corp., Security Response 5, 6 (2011), 29.
- [35] Mahsa Foruhandeh, Yanmao Man, Ryan Gerdes, Ming Li, and Thidapat Chantem. 2019. SIMPLE: Single-Frame Based Physical Layer Identification for Intrusion Detection and Prevention on in-Vehicle Networks. In Proceedings of the 35th Annual Computer Security Applications Conference (San Juan, Puerto Rico, USA) (ACSAC '19). Association for Computing Machinery, New York, NY, USA, 229–244. https://doi.org/10.1145/3359789.3359834
- [36] Aurélien Francillon and Claude Castelluccia. 2008. Code injection attacks on harvard-architecture devices. In Proceedings of the 15th ACM conference on Computer and communications security. 15–26.
- [37] Tommaso Frassetto, David Gens, Christopher Liebchen, and Ahmad-Reza Sadeghi. 2017. Jitguard: hardening just-in-time compilers with sgx. In Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security. 2405–2419.
- [38] Jiri Gaisler et al. 2001. The LEON processor user's manual. Gaisler research 2 (2001).
- [39] Shay Gal-On and Markus Levy. 2012. Exploring coremark a benchmark maximizing simplicity and efficacy. The Embedded Microprocessor Benchmark Consortium (2012).
- [40] Thanassis Giannetsos, Tassos Dimitriou, Ioannis Krontiris, and Neeli R Prasad. 2010. Arbitrary code injection through self-propagating worms in von neumann architecture devices. *Comput.* J. 53, 10 (2010), 1576–1593.
- [41] Zonghua Gu, Chao Wang, Ming Zhang, and Zhaohui Wu. 2014. WCET-aware partial control-flow checking for resource-constrained real-time embedded systems. IEEE Transactions on Industrial Electronics 61 (2014), 5652–5661. Issue 10. https://doi.org/10.1109/TIE.2014.2301752
- [42] Matthew R Guthaus, Jeffrey S Ringenberg, Dan Ernst, Todd M Austin, Trevor Mudge, and Richard B Brown. 2001. MiBench: A free, commercially representative embedded benchmark suite. In Proceedings of the fourth annual IEEE international workshop on workload characterization. WWC-4 (Cat. No. 01EX538). IEEE, 3–14.
- [43] Xiaochen Hao, Mingsong Lv, Jiesheng Zheng, Zhengkui Zhang, and Wang Yi. 2019. Integrating cyber-attack defense techniques into real-time cyber-physical systems. In 2019 IEEE 37th International Conference on Computer Design (ICCD). IEEE, 237–245.
- [44] Monowar Hasan, Sibin Mohan, Rodolfo Pellizzoni, and Rakesh B Bobba. 2017. Contego: An adaptive framework for integrating security tasks in real-time systems. arXiv preprint arXiv:1705.00138 (2017).
- [45] Wenjian He, Sanjeev Das, Wei Zhang, and Yang Liu. 2020. BBB-CFI: Lightweight CFI Approach Against Code-Reuse Attacks Using Basic Block Information. ACM Transactions on Embedded Computing Systems (TECS) 19, 1 (2020), 1–22.
- [46] ARM Holdings. [n.d.]. ARMv7-M Architecture Reference Manual, December 2014. https://developer.arm.com/documentation/ddi0403/latest/
- [47] Hong Hu, Shweta Shinde, Sendroiu Adrian, Zheng Leong Chua, Prateek Saxena, and Zhenkai Liang. 2016. Data-oriented programming: On the expressiveness of non-control data attacks. In 2016 IEEE Symposium on Security and Privacy (SP). IEEE, 969–986.
- [48] Wei Hu, Jason Hiser, Dan Williams, Adrian Filipi, Jack W Davidson, David Evans, John C Knight, Anh Nguyen-Tuong, and Jonathan Rowanhill. 2006.
 Secure and practical defense against code-injection attacks using software dynamic translation. In Proceedings of the 2nd international conference on Virtual execution environments. 2–12.
- [49] Zhijun Huang, Tao Zheng, Yunxiu Shi, and Ang Li. 2012. A dynamic detection method against ROP and JOP. In 2012 International Conference on Systems and Informatics (ICSAI2012). IEEE, 1072–1077.
- [50] Tomoaki Kawada, Shinya Honda, Yutaka Matsubara, and Hiroaki Takada. 2020. TZmCFI: RTOS-Aware Control-Flow Integrity Using TrustZone for Armv8-M. International Journal of Parallel Programming (2020), 1–21.
- [51] Sreenath Krishnadas. 2016. Concept and Implementation of AUTOSAR compliant Automotive Ethernet stack on Infineon Aurix Tricore board. (2016).
- [52] Donghyun Kwon, Jangseop Shin, Giyeol Kim, Byoungyoung Lee, Yeongpil Cho, and Yunheung Paek. 2019. uXOM: Efficient eXecute-Only Memory on {ARM} Cortex-M. In 28th {USENIX} Security Symposium ({USENIX} Security 19). 231–247.

- [53] Ruby B Lee, David K Karig, John P McGregor, and Zhijie Shi. 2004. Enlisting hardware architecture to thwart malicious code injection. In Security in Pervasive Computing. Springer, 237–252.
- [54] Jinfeng Li, Liwei Chen, Qizhen Xu, Linan Tian, Gang Shi, Kai Chen, and Dan Meng. 2020. Zipper stack: Shadow stacks without shadow. In European Symposium on Research in Computer Security. Springer. 338–358.
- [55] Hans Liljestrand, Thomas Nyman, Lachlan J. Gunn, Jan-Erik Ekberg, and N. Asokan. 2020. PACStack: an Authenticated Call Stack. arXiv:1905.10242 [cs.CR]
- [56] Arm Ltd. 2020. ARM Cortex-R. https://developer.arm.com/ip-products/processors/cortex-r
- [57] Microsoft Ltd. [n.d.]. Control Flow Guard Win32 apps. https://docs.microsoft.com/en-us/windows/win32/secbp/control-flow-guard
- [58] Minzhao Lyu, Dainel Sherratt, Arunan Sivanathan, Hassan Habibi Gharakheili, Adam Radford, and Vijay Sivaraman. 2017. Quantifying the Reflective DDoS Attack Capability of Household IoT Devices. In Proceedings of the 10th ACM Conference on Security and Privacy in Wireless and Mobile Networks (Boston, Massachusetts) (WiSec '17). Association for Computing Machinery, New York, NY, USA, 46-51. https://doi.org/10.1145/3098243.3098264
- [59] Stephen McCamant and Greg Morrisett. 2005. Efficient, verifiable binary sandboxing for a CISC architecture. (2005).
- [60] Mitra Nasri, Thidapat Chantem, Gedare Bloom, and Ryan M Gerdes. 2019. On the pitfalls and vulnerabilities of schedule randomization against schedule-based attacks. In 2019 IEEE Real-Time and Embedded Technology and Applications Symposium (RTAS). IEEE, 103–116.
- [61] Ahmad MK Nasser and Di Ma. 2020. SecMonQ: An HSM based security monitoring approach for protecting AUTOSAR safety-critical systems. Vehicular Communications 21 (2020), 100201.
- [62] Thomas Nyman, Jan-Erik Ekberg, Lucas Davi, and N Asokan. 2017. CFI CaRE: Hardware-supported call and return enforcement for commercial microcontrollers. In International Symposium on Research in Attacks, Intrusions, and Defenses. Springer, 259–284.
- [63] Nahmsuk Oh, Philip P Shirvani, and Edward J McCluskey. 2002. Control-flow checking by software signatures. IEEE transactions on Reliability 51, 1 (2002), 111–122.
- [64] James Pallister, Simon Hollis, and Jeremy Bennett. 2013. BEEBS: Open Benchmarks for Energy Measurements on Embedded Platforms. arXiv (2013), arXiv-1308.
- [65] Sandro Pinto and Nuno Santos. 2019. Demystifying arm trustzone: A comprehensive survey. ACM Computing Surveys (CSUR) 51, 6 (2019), 1–36.
- [66] Julien Proy, Karine Heydemann, Alexandre Berzati, and Albert Cohen. 2017. Compiler-Assisted Loop Hardening Against Fault Attacks. ACM Trans. Archit. Code Optim. 14, 4, Article 36 (Dec. 2017), 25 pages. https://doi.org/10.1145/3141234
- [67] Donald Ray and Jay Ligatti. 2012. Defining code-injection attacks. Acm Sigplan Notices 47, 1 (2012), 179–190.
- [68] Abhishek Rhisheekesan, Reiley Jeyapaul, and Aviral Shrivastava. 2019. Control flow checking or not? (for Soft Errors). ACM Transactions on Embedded Computing Systems 18 (2 2019). Issue 1. https://doi.org/10.1145/3301311
- [69] Gerardo Richarte et al. 2002. Four different tricks to bypass stackshield and stackguard protection. World Wide Web 1 (2002).
- [70] Ryan Roemer, Erik Buchanan, Hovav Shacham, and Stefan Savage. 2012. Return-oriented programming: Systems, languages, and applications. ACM Transactions on Information and System Security (TISSEC) 15, 1 (2012), 1–34.
- [71] Monireh Safari and Reihaneh Khorsand. 2018. PL-DVFS: combining Power-aware List-based scheduling algorithm with DVFS technique for real-time tasks in Cloud Computing. The Journal of Supercomputing 74, 10 (2018), 5578–5600.
- [72] Sonal Saha and Binoy Ravindran. 2012. An experimental evaluation of real-time DVFS scheduling algorithms. In Proceedings of the 5th Annual International Systems and Storage Conference. 1–12.
- [73] Dr Sarwar Sayeed, Hector Marco-Gisbert, Ismael Ripoll, and Miriam Birch. 2019. Control-Flow Integrity: Attacks and Protections. Applied Sciences (2019).
- [74] Simon Schuster, Peter Ulbrich, Isabella Stilkerich, Christian Dietrich, and Wolfgang Schröder-Preikschat. 2017. Demystifying soft-error mitigation by control-flow checking - A new perspective on its effectiveness. ACM Transactions on Embedded Computing Systems 16. Issue 5s. https://doi.org/10.1145/3126503
- [75] Hovav Shacham. 2007. 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*. 552–561.
- [76] Z. Shao, Q. Zhuge, Y. He, and E. H. . Sha. 2003. Defending embedded systems against buffer overflow via hardware/software. In 19th Annual Computer Security Applications Conference, 2003. Proceedings. 352–361. https://doi.org/10.1109/CSAC.2003.1254340
- [77] Amitabh Srivastava, Andrew Edwards, and Hoi Vo. 2001. Vulcan: Binary Transformation In A Distributed Environment. Technical Report MSR-TR-2001-50. 12 pages. https://www.microsoft.com/en-us/research/publication/vulcan-binary-transformation-in-a-distributed-environment/
- [78] Bowen Tang, Huan Ying, Wei Wang, and Huabin Tang. 2017. Eternal War in Software Security: A Survey of Control Flow Protection.
- [79] Victor Van der Veen, Dennis Andriesse, Enes Göktaş, Ben Gras, Lionel Sambuc, Asia Slowinska, Herbert Bos, and Cristiano Giuffrida. 2015. Practical context-sensitive CFI. In Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security.
- [80] Stijn Volckaert, Bart Coppens, and Bjorn De Sutter. 2015. Cloning your gadgets: Complete ROP attack immunity with multi-variant execution. IEEE Transactions on Dependable and Secure Computing 13, 4 (2015), 437–450.
- [81] Robert Wahbe, Steven Lucco, Thomas E Anderson, and Susan L Graham. 1993. Efficient software-based fault isolation. In *Proceedings of the fourteenth ACM symposium on Operating systems principles*. 203–216.
- [82] Robert J Walls, Nicholas F Brown, Thomas Le Baron, Craig A Shue, Hamed Okhravi, and Bryan C Ward. 2019. Control-flow integrity for real-time embedded systems. In 31st Euromicro Conference on Real-Time Systems (ECRTS 2019). Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik.

- [83] Tielei Wang, Tao Wei, Guofei Gu, and Wei Zou. 2010. TaintScope: A checksum-aware directed fuzzing tool for automatic software vulnerability detection. In 2010 IEEE Symposium on Security and Privacy. IEEE, 497–512.
- [84] Nathanael R Weidler, Dane Brown, Samuel A Mitchell, Joel Anderson, Jonathan R Williams, Austin Costley, Chase Kunz, Christopher Wilkinson, Remy Wehbe, and Ryan Gerdes. 2019. Return-oriented programming on a resource constrained device. Sustainable Computing: Informatics and Systems 22 (2019), 244–256.
- [85] Penglin Yang, Limin Tao, and Haitao Wang. 2018. RTTV: a dynamic CFI measurement tool based on TPM. IET Information Security 12, 5 (2018), 438–444
- [86] Man-Ki Yoon, Sibin Mohan, Chien-Ying Chen, and Lui Sha. 2016. Taskshuffler: A schedule randomization protocol for obfuscation against timing inference attacks in real-time systems. In 2016 IEEE Real-Time and Embedded Technology and Applications Symposium (RTAS). IEEE, 1–12.
- [87] Ming Zhang, Zonghua Gu, Hong Li, and Nenggan Zheng. 2018. WCET-Aware Control Flow Checking with Super-Nodes for Resource-Constrained Embedded Systems. IEEE Access 6 (7 2018), 42394–42406. https://doi.org/10.1109/ACCESS.2018.2852805
- [88] Mingwei Zhang and R Sekar. 2013. Control flow integrity for {COTS} binaries. In 22nd {USENIX} Security Symposium ({USENIX} Security 13). 337–352.
- [89] Jie Zhou, Yufei Du, Zhuojia Shen, Lele Ma, John Criswell, and Robert J. Walls. 2020. Silhouette: Efficient Protected Shadow Stacks for Embedded Systems. In 29th USENIX Security Symposium (USENIX Security 20). USENIX Association, 1219–1236. https://www.usenix.org/conference/usenixsecurity20/presentation/zhou-jie
- [90] Nikola Zlatanov. 2016. ARM Architecture and RISC Applications. (2016).