

# High-Performance Software Protection using Reconfigurable Architectures

Joseph Zambreno, *Student Member, IEEE*, Dan Honbo, *Student Member, IEEE*,  
 Alok Choudhary, *Senior Member, IEEE*, Rahul Simha, *Member, IEEE*, and Bhagi Narahari, *Member, IEEE*

**Abstract**—One of the key problems facing the computer industry today involves ensuring the integrity of end-user applications and data. Researchers in the relatively new field of software protection investigate the development and evaluation of controls that prevent the unauthorized modification or use of system software. While many previously developed protection schemes have provided a strong level of security, their overall effectiveness has been hindered by a lack of transparency to the user in terms of performance overhead. Other approaches take to the opposite extreme and sacrifice security for the sake of this need for transparency. In this work we present an architecture for software protection that provides for a high level of both security and user transparency by utilizing Field Programmable Gate Array (FPGA) technology as the main protection mechanism. We demonstrate that by relying on FPGA technology, this approach can accelerate the execution of programs in a cryptographic environment, while maintaining the flexibility through reprogramming to carry out any compiler-driven protections that may be application-specific. Furthermore, we show how programmable FPGA resources not reserved towards software protection can be realized as performance-oriented architectural optimizations, and we evaluate the effectiveness of this concept with an investigation into instruction prefetching.

## I. INTRODUCTION AND MOTIVATION

**T**HREATS to a particular piece of software can originate from a variety of sources. A substantial problem from an economic perspective is the unauthorized copying and redistribution of applications, otherwise known as *software piracy*. Although the actual damage sustained from the piracy of software is certainly a debatable matter, some industry watchdog groups have estimated that software firms in 2002 lost as much as \$2 billion in North American sales alone [1]. A threat that presents a much more direct harm to the end-user is *software tampering*, whereby a hacker maliciously modifies and redistributes code in order to cause large-scale disruptions in software systems or to gain access to critical information. For these reasons, *software protection* is considered one of the more important unresolved research challenges in security today [2]. In general, any software protection infrastructure should include 1) a method of limiting an attacker’s ability to understand the higher level semantics of an application given a low-level (usually binary) representation, and 2) a system

J. Zambreno, D. Honbo, and A. Choudhary are with the Department of Electrical and Computer Engineering, Northwestern University, Evanston, IL 60208 USA (e-mail: zambro1@ece.northwestern.edu; d-honbo@northwestern.edu; choudhar@ece.northwestern.edu).

R. Simha and B. Narahari are with the Department of Computer Science, The George Washington University, Washington, DC 20052 USA (e-mail: simha@gwu.edu; narahari@gwu.edu).

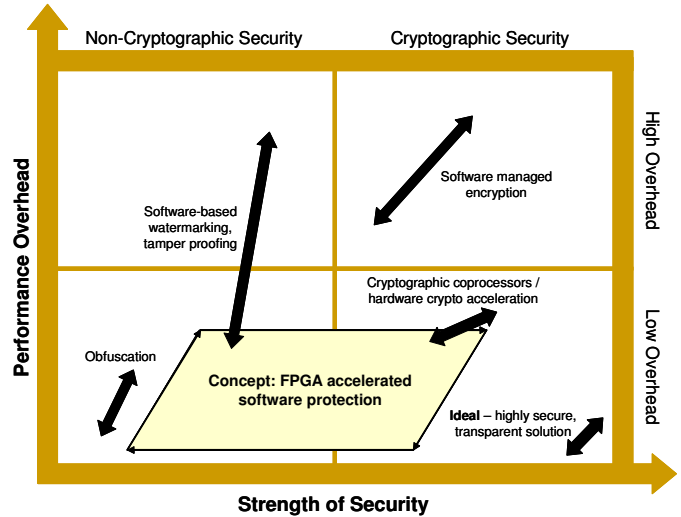


Fig. 1. Performance and security strength of various software protection approaches

of checks that make it suitably difficult to modify the code at that low level. When used in combination, these two features can be extremely effective in preventing the circumvention of software authorization mechanisms.

Current approaches to software protection can be categorized both by the strength of security provided and the performance overhead when compared to an unprotected environment (see Fig. 1). Two distinct categories emerge from this depiction - on one end of the security spectrum are the solely compiler-based techniques that implement both static and dynamic code validation through the insertion of objects into the generated executable. On the other end are the somewhat more radical methods that encrypt all instructions and data and that often require the processor to be architected with cryptographic hardware. Both of these methods have some practical limitations. The software-based techniques will effectively only hinder an attacker, since tools can be built to identify and circumvent the protective checks. On the other hand, the cryptographic hardware approaches, while inherently more secure, are limited in the sense that their practical implementation requires a wholesale commitment to a custom processor technology. More background on these software protection schemes can be found in Section 2.

Also, as can be inferred from Fig. 1, software protection is not an all-or-nothing concept, as opposed to other aspects

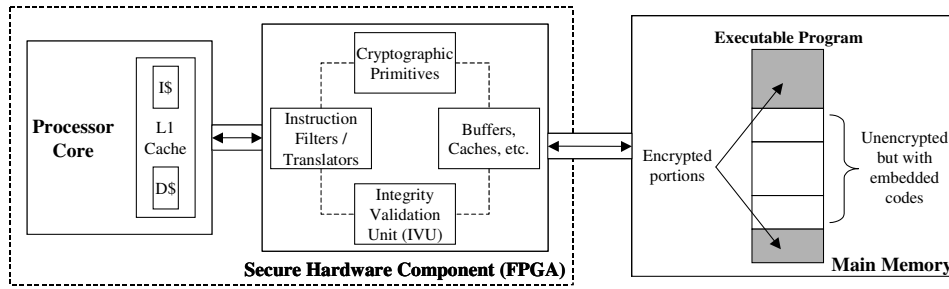


Fig. 2. Conceptual view

of computer security where the effectiveness of an approach depends on its mathematical intractability; the resultant performance overhead is often not a key design consideration of these systems. In contrast, performance is equally important to the strength of security when protecting user applications, as there is generally some level of user control over the system as a whole. Consequently, any software protection scheme that is burdensome from a performance perspective will likely be turned off or routed around.

Field Programmable Gate Arrays (FPGAs) are a hardware resource that combine various amounts of user-defined digital logic with customizable interconnect and I/O. A key feature of FPGAs is that their functionality can be reconfigured on multiple occasions, allowing for changes in the design to be implemented after the initial time of development. FPGAs have become an increasingly popular choice among architects implementing algorithms in fields such as multimedia processing or cryptography - this has been attributed to the fact that the design process is much more streamlined than that for ASICs, as FPGAs are a fixed hardware target.

In this paper we present a high-performance architecture for software protection that uses this type of reconfigurable technology. By choosing FPGAs as the main protection mechanism, this approach is able to merge the application tunability of the compiler-based approaches with the additional security that comes with a hardware implementation. As can be seen in Fig. 2, our proposed method works by supplementing a standard processor with an FPGA-based Integrity Validation Unit (IVU) that sits between the highest level of on-chip cache and main memory. This IVU is capable of performing fast decryption similar to any other hardware accelerated cryptographic coprocessing scheme, but more importantly the IVU also has the ability to recognize and certify binary messages hidden inside regular unencrypted instructions. Consequently our approach is completely complementary to current code restructuring techniques found in the software protection literature, with the added benefit that as the run-time code checking can be performed entirely in hardware it will be considerably more efficient. Our experiments show, that for most of our benchmarks, the inclusion of the FPGA within the instruction stream incurs a performance penalty of less than 20%, and that this number can be greatly improved upon with the utilization of unreserved reconfigurable resources for architectural optimizations, such as buffering and prefetching.

The remainder of this paper is organized as follows. In

Section 2 we provide additional background into the field of software protection, with a review of some of the more commonly-used defensive techniques. Section 3 sets the backdrop of our research, both illustrating the threat model under which we apply our approach and making the argument for the level of security that is provided. In Section 4 we present our architecture in more detail, illustrating how we can utilize an FPGA situated in the instruction stream to ensure software integrity. In this section we also provide an introduction to some custom compiler techniques that are well suited to such an architecture. Section 5 discusses the performance implication of our approach, first by explicitly quantifying the security/performance tradeoff and then through experimental analysis. In Section 6, we describe several architectural optimizations that are possible within our framework and provide results detailing the effectiveness of instruction prefetching. Finally, in Section 7 we present our conclusions alongside a discussion of future techniques that are currently in development.

## II. BACKGROUND

While the term *software protection* often refers to the purely software-based defense mechanisms available to an application after all other safeguards have been broken, in practical systems this characterization is not necessarily so precise. For our purposes, we classify hardware-supported secure systems as being tools for software protection as long as they include one of the three commonly found elements of a software-based protection scheme [3] detailed below. A survey of the broader area of software security in the context of DRM appears in [4], [5].

**Watermarking** is a technique whereby messages are hidden inside a piece of software in such a way that they can be reliably identified [6]. While the oldest type of watermarking is the inclusion of copyright notices into both code and digital media, more recent watermarking approaches have focused on embedding data structures into an application, the existence of which can then be verified at run-time. Venkatesan et al. present an interesting variation on watermarking in [7], where the watermark is a subprogram that has its control flow graph merged with the original program in a stealthy manner.

The concept behind **tamper-proofing** is that a properly secured application should be able to safely detect at run-time if it has been altered. A type of dynamic “self-checking” is proposed in both [8] and [9], these approaches assert

application integrity by essentially inserting instructions to perform code checksums during program execution. An interesting technique is proposed by Aucsmith in [10], in which partitioned code segments are encrypted and are handled in a fashion such that only a single segment is ever decrypted at a time. The novelty of this approach is that the program flows between segments based on a function of the current segment; consequently any modification to the code will result in an incorrect path execution. One flaw with many of these tamper-proofing approaches is that in most architectures it is relatively easy to build an automated tool to reveal the checking mechanisms. For example, checksum computations can be easily identified by finding code that operates directly on the instruction address space. Accordingly, the security of these approaches is depends heavily on the security of the checking mechanisms themselves.

These same checksums suggested for software can be computed in hardware. Unfortunately, this approach relies on designing hardware that has access to memory and that is started by software instructions. The latter creates an opening for attack whereas the former complicates both design and efficiency since the specialized hardware must now have access to the memory address and data lines.

Proof-Carrying Code (PCC) is a recently proposed solution that has techniques in common with other tamper-proofing approaches. PCC allows a host to verify code from an untrusted source [11], [12], [13], [14]. Safety rules, as part of a theorem-proving technique, are used on the host to guarantee proper program behavior. Applications include browser code (applets) [15] and even operating systems [14]. One advantage of proof-carrying software is that the programs are self-certifying, independent of encryption or obscurity. The PCC method is essentially a self-checking mechanism and is vulnerable to the same problems that arise with the code checksum methods discussed earlier; in addition they are static methods and do not address changes to the code after instantiation.

The goal of **obfuscation** is to limit code understanding through the deliberate mangling of program structure - a survey of such techniques can be found in [16]. Obfuscation techniques range from simple encoding of constants to more complex methods that completely restructure code while maintaining correctness [16], [17]. Wang et al. [18], [19] provide some transformations that make it difficult to determine the control flow graph of the original program, and show that determining the control flow graph of the transformed code is NP-hard. It is important to note that, just as with tamper-proofing, obfuscation can only it make the job of an attacker more difficult, since tools can be built to automatically look for obfuscations, and tracing through an executable in a debugger can reveal vulnerabilities. These and other theoretical limitations are discussed in more detail in [20].

#### A. Other Hardware-based Approaches

Using our definition, there have been several hardware-based software protection approaches. **Secure coprocessors** are computational devices that enable execution of encrypted programs. Programs, or parts of the program, can be run in

an encrypted form on these devices thus never revealing the code in the untrusted memory and thereby providing a tamper resistant execution environment for that portion of the code. A number of secure coprocessing solutions have been designed and proposed, including systems such as IBM's Citadel [21], Dyad [22], [23], [24], the Abyss and mAbyss systems [25], [26], [27], and the commercially available IBM 4758 which meets the FIPS 140-1 Level 4 validation [28], [29], [30].

**Smart cards** can also be viewed as type of secure coprocessing; a number of studies have analyzed the use of smart cards for secure applications [31], [32]. Sensitive computations and data can be stored in the smart card but they offer no direct I/O to the user. Most smart card applications focus on the secure storage of data although studies have been conducted on using smart cards to secure an operating system [33]. As noted in [8], smart cards can only be used to protect small fragments of code and data.

Recent commercial hardware security initiatives have focused primarily on cryptographic acceleration, domain separation, and trusted computing. These initiatives are intended to protect valuable data against software-based attacks and generally do not provide protection against physical attacks on the platform. MIPS and VIA have added cryptographic acceleration hardware to their architectures. MIPS Technologies' SmartMIPS ASE [34] implements specialized processor instructions designed to accelerate software cryptography algorithms, while VIA's Padlock Hardware Security Suite [35] adds a full AES encryption engine to the processor die. Both extensions seek to eliminate the need for cryptographic coprocessors.

Intel's LaGrande Technology [36], ARM's TrustZone Technology [37], and MIPS Technologies' SmartMIPS ASE implement secure memory restrictions in order to enforce domain separation. These restrictions segregate memory into secure and normal partitions and prevent the leakage of secure memory contents to foreign processes. Intel and ARM further strengthen their products' domain separation capabilities by adding a processor privilege level. A Security Monitor process is allowed to run at the added privilege level and oversee security-sensitive operations. The Security Monitor resides in protected memory and is not susceptible to observation by user applications or even the Operating System.

Several companies have formed the so-called Trusted Computing Group (TCG) to provide hardware-software solutions for software protection [38]. The TCG defines specifications for the Trusted Platform Module, a hardware component that provides digital signature and key management functions, as well as shielded registers for platform attestation. Intel's LaGrande Technology and Microsoft's Next-Generation Secure Computing Base [39] combine the TPM module with processor and chipset enhancements to enable platform attestation, sealed storage, strong process isolation, and secure I/O channels. All of these approaches require processor or board manufacturers to commit to a particular design, and once committed, are locked into the performance permitted by the design.

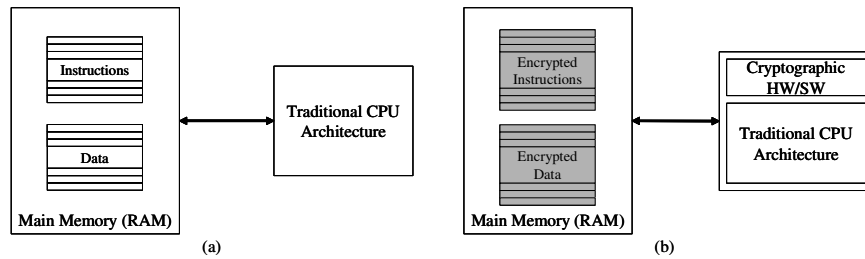


Fig. 3. Standard von Neumann architecture with a CPU and main memory, as compared to (b) - an Encrypted Execution and Data (EED) platform

### B. Closely Related Work

In [40], researchers at Stanford University proposed an architecture for tamper-resistant software based on an eXecute-Only Memory (XOM) model that allows instructions stored in memory to be executed but not manipulated. A hardware implementation is provided that is not dissimilar to our proposed architecture, with specialized hardware being used to accelerate cryptographic functionality needed to protect data and instructions on a per-process basis. Three key factors differentiate our work from the XOM approach. One distinction is that our architecture requires no changes to the processor itself. Also, our choice of reconfigurable hardware permits a wide range of optimizations that can shape the system security and resultant performance on a per-application basis. Most importantly, we consider a host of new problems arising from attacks on encrypted execution and data platforms.

In [41], researchers at UCLA and Microsoft Research propose an intrusion prevention system known as the Secure Program Execution Framework (SPEF). Similar to our proposed work, the SPEF system is used as the basis for compiler transformations that both obfuscate and also embed integrity checks into the original application that are meant to be verified at run-time by custom hardware. While an interesting project, the SPEF work in its current form concentrates solely on preventing intruder code from being executed, and consequently neglects similar attacks that would focus mainly on data integrity. Also, the compiler-embedded constraints in the SPEF system require a predefined hardware platform on which to execute; this limits the scope of any such techniques to the original processor created for such a purpose.

Pande et al. [42], [43] address the problem of information leakage on the address bus wherein the attacker would be snooping the address values to gain information about the control flow of the program. They provide a hardware obfuscation technique which is based on dynamically randomizing the instruction addresses. This is achieved through a secure hardware coprocessor which randomizes the addresses of the instruction blocks, and rewrites them into new locations. Using this scheme, for example, the attacker would be unable to determine control flows such as loops and branches since they do not see a recurrence in the addresses being fetched. While this scheme provides obfuscation of the instruction addresses, thereby providing a level of protection against IP theft, it does not prevent an attacker from injecting their own instructions to be executed and thereby disrupting the processor and application.

### C. FPGAs and Security

FPGAs have been used for security-related purposes in the past as hardware accelerators for cryptographic algorithms. Along these lines Dandalis and Prasanna [44], [45] have led the way in developing FPGA-based architectures for internet security protocols. Several similar ideas have been proposed in [46], [47]. The FPGA manufacturer Actel [48] offers commercial IP cores for implementations of the DES, 3DES, and AES cryptographic algorithms. Similar to our approach, these implementations utilize FPGAs not only for their computational speed but for their programmability; in security applications the ability to modify algorithmic functionality in the field is crucial. In an attempt to raise consumer confidence in FPGA security, Actel is currently developing new anti-fuse technologies that would make FPGAs more difficult to reverse-engineer [49].

## III. SYSTEM MODEL AND RATIONALE

Fig. 3a shows a depiction of the standard von Neumann architecture with a CPU and main memory (RAM); we will refer to the main memory as RAM in the sequel. In the standard model, a program consisting of instructions and data is placed in RAM by a loader or operating system. Then, the CPU fetches instructions and executes them. Apart from the complexity introduced by multiprocessors and threads, this basic model applies to almost any computing system today including desktops, servers and small embedded devices.

However, from a security point of view, the execution of an application is far from safe. An attacker with access to the machine can not only examine the program (information leakage) but can actively interfere with the execution (disruption) while also accumulating information useful in attacks on similar systems.

Most commercial desktop and server systems are today being designed to sustain attacks via a network. The bulk of these attacks are from hackers or commercial predators who are not on-site. However, military systems need to be designed to also sustain the sophisticated attacks possible when the computing equipment is captured and carefully taken apart in a resourceful laboratory by experts. For such systems, it is desirable to encrypt the application to minimize information leakage and to prevent disruption.

The starting point for our proposed work is a recent body of work that has proposed building computing platforms with encrypted execution. We refer to these as EED (Encrypted Execution and Data) platforms. In this platform (Fig. 3b),

the executable and application are encrypted. The processor (or supporting hardware) is assumed to have a key. These processors are either separate chips or processor cores that are supplied by core manufacturers for use in System-on-Chip (SoC) systems. The overall goal is to prevent leakage of information, to prevent tampering and to prevent disruption. Consider, for example, the computing platform of a wireless device. The chipset will typically consist of a processor chip, a transceiver, RAM, and supporting connector chips and microcontrollers. We assume that an attacker that gains control of the device has significant hardware resources available for reverse engineering and malicious insertion of signals between chips. For the highest degree of security, both instructions and data will need to be encrypted using well-established encryption techniques. It should be noted that full-fledged EED platforms are still in their infancy.

The basis for our proposed work is the following:

- EED platforms, while undoubtedly more secure than the standard von Neumann model, are nonetheless still vulnerable to attacks that do not need decryption. That is, the attacker can find vulnerabilities without needing to decrypt and understand the software. We will refer to these attacks as EED attacks.
- Because attackers are presumed to be sophisticated, neither the RAM nor the CPU can be fully trusted. This situation also arises when RAM and CPU manufacturing is sub-contracted to third parties whose designs or cores cannot be easily verified.
- A class of reconfigurable hardware technology, the Field Programmable Gate Array (FPGA) has proved adept at solving performance-related problems in many computing platforms. As a result, tested FPGAs are commercially available for a variety of processors. Our approach exploits the combination of the programmability of FPGAs with the inherent additional security involved with computing directly in hardware to address EED attacks.
- A key part of our exploiting FPGAs involves the use of compiler technology to analyze program structure to enable best use of the FPGA, to help address key management and to increase performance. Consequently our approach allows for a level of security that is tunable to an individual application.

In a nutshell, the addition of the FPGA to the von Neumann model results in a new platform to sustain EED attacks.

**A. Focus Areas**

The contributions of our work can be categorized into four areas as seen in Fig. 4. We use the term *structural integrity* to refer to the proper execution path of a program when the data is assumed to be correct. Since an EED attacker can alter the control flow without decryption or even touching the data, we refer to such an attack as a structural EED attack. Thus the first area of contributions is shown as Area 1 in the figure, in which we use a compiler-FPGA approach to address structural attacks.

The second area of contribution arises from considering EED attacks on *data integrity*. There are two sub-categories

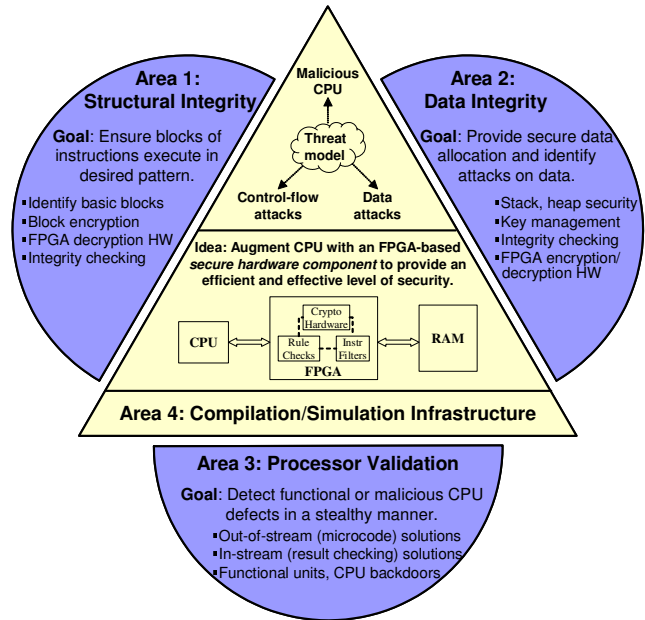


Fig. 4. Research focus areas as they relate to the conceptual view

here, the regular data used by an application and the runtime data (stack, heap) needed for the execution. Thus, our contribution in this second area is the use of the compiler-FPGA approach to provide key management techniques for data and data integrity techniques for runtime data.

The third area, *processor validation*, arises from considering a maliciously inserted processor instrumented to allow the attacker control over functional operations. Our approach is to have the compiler instrument the code with processor-validation meta-instructions that are then interpreted in the FPGA to validate desired processor operations. Finally, perhaps most importantly, we developed a compiler-FPGA infrastructure for the purpose of implementing and testing our ideas. This infrastructure, which features a modern processor (ARM family) and compiler (gcc), can be used beyond our project.

**B. Threat Model**

As mentioned in Section 1, research in the field of software protection has seen its motivation arrive from two different directions. On one side are vendors of various types of electronic media - their main concern is the unauthorized use or copying of their product. On the other side are those end users (with supporting hardware and software vendors) whose main interest is in protecting personal and corporate systems from outside attack. While the goals may be different, the approaches used by hackers in avoiding Digital Rights Management (DRM) schemes are often quite similar to those used by malicious crackers in attacking web servers and other unsecured applications. Hackers around the world know that the first step in attacking a software system is to first understand the software through the use of a debugger or other tracing utilities, and then to tamper with the software to enable a variety of exploits. Common means of exploiting software

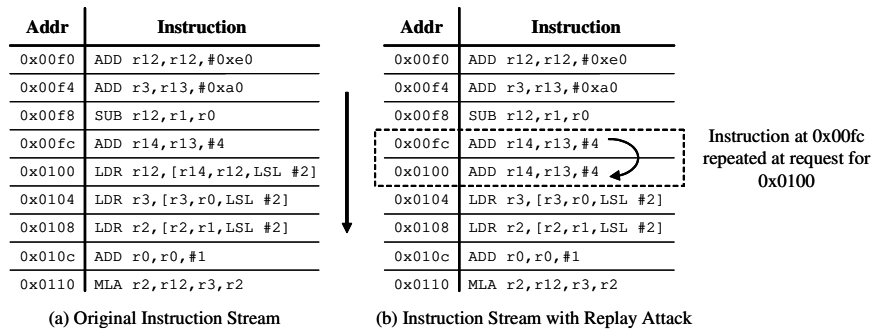


Fig. 5. A sample instruction stream that is targeted with a replay attack

include buffer overflows, malformed printf() statements, and macro viruses.

Consider a sophisticated attacker who has gained control of an EED computing platform. For example, this platform could be a CPU board on a desktop or a board from an embedded device. In a resourceful laboratory, the attacker can control the various data, address and control lines on the board and will have access to the designs of well-known commercial chips. Using this information, the attacker can actively interfere in the handshaking protocol between CPU and RAM, can insert data into RAM or even switch between the actual RAM and an attacker’s RAM during execution.

We will assume that the cryptographic strength is such that the attacker cannot actually decrypt the executable or data. How then does an attacker disrupt a system without being able to decrypt and understand? The following are examples of EED attacks:

- **Replay attacks** - consider an EED platform with encrypted instructions. An attacker can simply re-issue an instruction from the RAM to the CPU. This is relatively easy for an attacker who has complete control of the processor-RAM bus. What does such an attack achieve? The application program logic can be disrupted and the resulting behavior exploited by the attacker. Indeed, by observing the results of a multitude of replay attacks, the attacker can catalogue information about the results of individual replay attacks and use such attacks in tandem for greater disruption. A sample replay attack is depicted in Fig. 5.
- **Control-flow attacks** - the sophisticated attacker can elucidate the control-flow structure of a program without decryption. This can be done by simply sniffing the bus, recording the pattern of accesses and extracting the control-flow graph from the list of accesses. To disrupt, the attacker can prematurely transfer control out of a loop, or even simply transfer control to a distant part of the executable. Again, by repeatedly studying the impact of such control-flow attacks, the attacker can accumulate a database of attacks to be used in concert. A sample control-flow attack is depicted in Fig. 6.
- **Runtime data attacks** - by examining the pattern of data write-backs to RAM, the attacker can intelligently guess the location of the run-time stack even if that is encrypted, as commonly used software stack structures are relatively

simple (Fig. 7). By swapping contents in the stack, the attacker can disrupt the flow of execution or parameter passing via the stack. Again, the attacker does not need to decrypt to achieve this disruption. Similarly, by observing the same pattern of accesses, the attacker can guess at heap locations and engage in similar attacks.

- **Application data attacks** - an attacker can simply change application data (using, for example, replays from the time-history of a variable) to create an attack. This might cause improper execution (i.e., change the flow of execution) or the computation of incorrect results.
- **Improper processor computations** - the CPU itself may be untrustworthy, since an attacker with considerable resources may simulate the entire CPU and selectively change outputs back to RAM.

Taken together, the above attacks can also be combined with cryptanalytic techniques to uncover cribs for decryption. This suggests that a secure computing platform should be able to detect such attacks and prevent disruption.

### C. A Note on FPGA Security

Thus far we have been working off the assumption that as the software protection mechanisms are situated directly in hardware (FPGA), these dynamic checks are inherently more secure than those validated through purely-software mechanisms. While this may very well be the case, some recent work has suggested that the use of FPGA hardware in embedded systems invites a range of attacks. For example, physical attacks that take advantage of the fact that not all data is lost in a Static RAM (SRAM) memory cell when its power is turned off have been shown to be very successful at retrieving information from memory devices that use SRAM technology [50]. It is likely that the techniques developed to facilitate these types of attacks would also be useful in extracting information from SRAM FPGAs. Another point of weakness is the FPGA bitstream. In our system, this file would need to be stored on disk just as the application itself. This potentially exposes some of our techniques to a reverse engineering attack.

Before we lose hope in the prospect of an FPGA as a secure hardware component, several clarifications need to be made. First, many of the proposed attacks on FPGAs assume that the physical security of the entire system has already been breached. No system, regardless of the underlying hardware

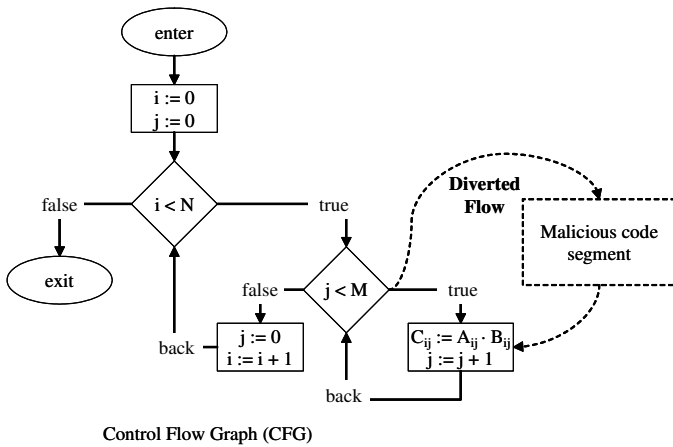


Fig. 6. A sample Control Flow Graph (CFG) with a diverted flow

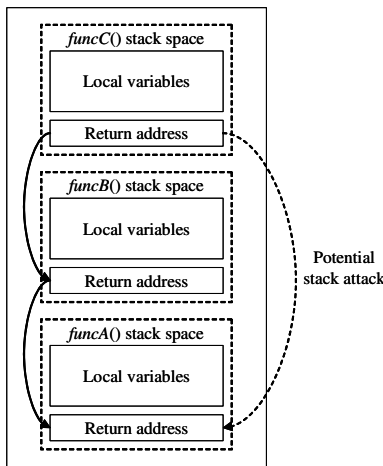


Fig. 7. A structural attack on an instruction stack

or software security, is safe assuming that the physical security has been breached. Consequently, under such assumptions, the choice of FPGA provokes no further risks in terms of physical security. Second, reverse engineering a hardware description is a considerably more difficult task than reverse-engineering an equivalent software description, due to the lack of readily available tools and the increased complexity of hardware designs. This task can be made even more difficult through the encrypting of the bitstream itself. These reasons, combined with the fact that newer FPGAs are being built with highly secure antifuse and flash-based technologies, have motivated the argument made by some that FPGAs are more difficult to reverse-engineer than an equivalently functional ASIC [48].

IV. OUR APPROACH

At its highest level, our approach is best classified as an attempt to combine the tunability of classical compiler-based software protection techniques with the additional security afforded through hardware. Consequently, our work is unique in that we utilize a combined hardware/software technique, and that we provide tremendous flexibility to application designers in terms of positioning on the security/performance spectrum. Also, going back to Fig. 1, our approach improves upon

previous software protection attempts by accelerating their performance while not sacrificing any security.

A. Architecture Overview

We accomplish our low-overhead software protection this through the placement of an FPGA between the highest level of on-chip cache and main memory in a standard processor instruction stream (see Fig. 8). In our architecture, the FPGA traps all instruction cache misses, and fetches the appropriate bytes directly from higher-level memory. These instructions would then be translated in some fashion and possibly verified before the FPGA satisfies the cache request. Both the translation and verification operations could be customized

The key features of our software protection architecture can be summarized as follows:

- **Fast Decryption** - Much work has been done in recent years on FPGA implementations of cryptographic algorithms; a popular target is the Advanced Encryption Standard (AES) candidate finalists [51]. These results have shown that current-generation FPGAs are capable of extremely fast decryption, as an example in [52] an implementation of Rijndael, the AES cipher [53] attained a throughput far exceeding 10 Gbps when fully pipelined. Consequently, by placing AES hardware on our FPGA we can perform instruction stream decryption without the prohibitive delays inherent in an equivalent software implementation.
- **Instruction Translation** - As a computationally less complex alternative to the full decryption described above, we can instantiate combinational logic on the FPGA to perform a binary translation of specific instruction fields. For example, a simple lookup table can map opcodes such that all addition instructions become subtractions and vice versa. As the mapping would be completely managed by the application developer, this technique would provide an effective yet extremely low-cost level of obfuscation.
- **Dynamic Validation** - In order to effectively prevent tampering, we can program the FPGA with a set of “rules” that describe correct program functionality. For a code segment that is fully encrypted, these rules can be as simple as requiring that each instruction decode properly. An example of a more complex rule would be one based on a predetermined instruction-based watermark. These software integrity checks can be performed either at each instruction or at a predefined interval - in either case given a detection of a failure condition the FPGA would require special attention from the system. Although it is possible that the FPGA would set a control signal to halt the CPU directly, it is more practical that it would simply return hard-coded instructions that can warn the operating system of a security violation.
- **Message Passing** - With the instantiation of simple instruction decode logic on the FPGA, we can utilize undefined mappings in the processor’s instruction set architecture to pass directives to the FPGA. Unlike most ISA modifications, this would require no changes in the

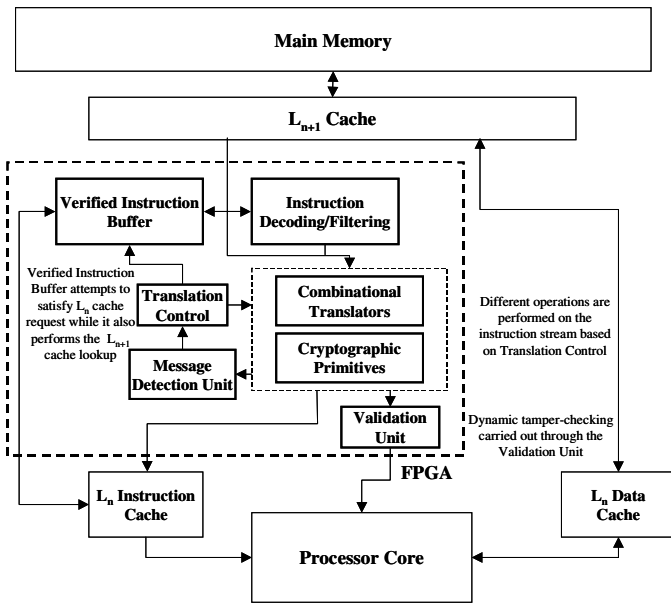


Fig. 8. FPGA-Based software protection architecture

processor itself, as the passed message can actually be returned as just a normal NOP instruction to the  $L_n$  cache. Some potentially useful messages include data-embedded watermark addresses, cryptographic keys, and directives for translation mode shifting.

- **Performance Optimizations** - Many of the operations required by these components are relatively simple in terms of required instantiated logic. Consequently, for any choice of FPGA technology it is quite possible that a portion of the total reconfigurable resources will remain unused. As our FPGA is situated in the instruction fetch stream, it makes sense to utilize those “extra” resources as performance-oriented architectural optimizations. In Section 6, we investigate into some depth the performance impact of implementing instruction prefetching [54], [55], [56].

### B. Compiler Support

As mentioned earlier, every feature that is implemented in our architecture requires extensive compiler support to both enhance the application binary with the desired code transformations and to generate a hardware configuration for an FPGA that incorporates the predetermined components. Many of the standard stages in a compiler can be modified to include software protection functionality. As an example, consider the data flow analysis module. In a standard performance-oriented compiler, the goal of this stage is to break the program flow into basic blocks and to order those blocks in a manner that increases instruction locality. However, it is possible to reorganize the code sequences in such a fashion that sacrifices some performance in exchange for an increased level of obfuscation. This is an example of a transformation that requires no run-time support; for others (such as those that use encryption), the compiler must also generate all the information needed to configure an FPGA that can reverse

the transformations and verify the output. Our work currently focuses on compiler back-end transformations (i.e. data flow analysis, register allocation, code generation), although there are possibilities in the front-end to examine in the future as well.

### C. Example Implementations

Up to this point, in detailing the different features available in our software protection architecture, we have maintained a fairly high-level view. In the following sections, we delve into the specifics of two examples that illustrate the potential of our combined compiler/FPGA technique.

**Example 1 - Tamper-resistant register encoding.** Consider a code segment in any application and focus on the instructions as they are fetched from memory to be executed by the processor (the *instruction stream*). In most instruction set architectures, the set of instructions comprising a sequentially-executed code segment will contain instructions that use registers. In this example, the decoding unit in the FPGA will extract one register in each register-based instruction - we can refer to the sequence of registers so used in the instruction stream as the *register stream*. In addition, the FPGA also extracts the opcode stream. As an example, the sequence of instructions *load x, r1 | load y, r2 | add r1, r2* can be decoded to generate the register stream of  $r_1, r_2, r_2$ . After being extracted, the register stream is then given a binary encoding; in our example  $r_1$  can encode **0** and  $r_2$  can encode **1**, and therefore the particular sequence of registers corresponds to the code **0 1 1**.

A binary translation component then feeds this string through some combinational logic to generate a key. This function can be as simple as just flipping the bits, although more complex transformations are possible. The key is then compared against a function of the opcode stream. In this example, an instruction filter module picks out a branch instruction following the register sequence and then compares the key to a hash of the instruction bits. If there is a match, the code segment is considered valid, otherwise the FPGA warns the system that a violation condition has been detected. This concept is similar to those proposed in [41] and [57].

How does such an approach work? As register-allocation is performed by the compiler, there is considerable freedom in selecting registers to allow for any key to be passed to the FPGA (the registers need not be used in contiguous instructions since it is only the sequence that matters). Also, the compiler determines which mechanism will be used to filter out the instructions that will be used for comparisons. If the code has been tampered with, there is a very high probability that the register sequence will be destroyed or that the opcode filtering will pick out a different instruction. As an example, if the filtering mechanism picks out the sixth opcode following the end of a register sequence of length  $k$ , any insertion or deletion of opcodes in that set of instructions would result in a failure.

It is important to note that this type of tamper-proofing would not normally be feasible if implemented entirely in software, since the checking computation could be easily identified



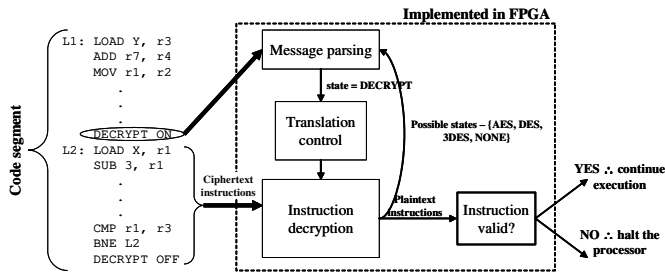


Fig. 9. Selective Basic Block Encryption

and avoided entirely. Also, the register sequence can be used to encode several different items, such as authorization codes or cryptographic keys. This technique can also be used to achieve code obfuscation by using a secret register-to-register mapping in the FPGA. Thus, if the FPGA sees the sequence  $(r_1, r_2, r_1)$ , this can be interpreted by the FPGA as an intention to actually use  $r_3$ . In this manner, the actual programmer intentions can be concealed by using a mapping customized to a particular processor. Finally, we note that when examining a block of code to be encrypted using our scheme, it will often be the case that the compiler will lack a sufficient number of register-based instruction with which to encode a suitable key. In this case, “place-holder” instructions which contain the desired sequence values, but which otherwise do not affect processor state, will need to be inserted by the compiler. In Section 5.2 we examine how these inserted instructions can effect the overall system performance.

**Example 2 - Selective basic block encryption.** Selective encryption is a useful technique in situations where certain code segments have high security requirements and a full-blown cryptographic solution is too expensive from a performance perspective. The example code segment in Fig. 9 shows how the compiler could insert message instructions to signify the start of an encrypted basic block (the compiler would also encrypt the block itself). As this message is decoded by the FPGA, an internal state would be set that directs future fetched instructions to be fed into the fast decryption unit. These control signals could be used to distinguish between different ciphers or key choices. The freshly-decrypted plaintext instructions would then be validated before being returned to the Ln cache of Fig. 8. The encryption mode could then be turned off or modified with the decoding of another message instruction.

Although properly encrypting a code segment makes it unreadable to an attacker who does not have access to the key, using cryptography by itself does not necessarily protect against tampering. A simple way of verifying that instructions have not been tampered with is to check if they decode properly based on the original instruction set architecture specification. However, this approach does not provide a general solution, as the overwhelming portion of binary permutations are usually reserved for the instruction set of most processors. This increases the likelihood that a tampered ciphertext instruction would also decode properly. A common approach is the use of one-way hash functions (the so-called “message digest” functions), but in our case, it would

be prohibitively slow to calculate the hash of every encrypted basic block in even medium-sized applications. A more simple approach would be to recognize patterns of instructions in the code segment that make sense in terms of the register access patterns. Specific care must also be taken to ensure the integrity of the message instructions themselves. This can be implemented through a combination of the register encoding techniques discussed previously and other dynamic code checking methods.

This example describes an approach that would be inherently slow in a purely-software implementation. Consequently, using the FPGA for decryption allows the application designer the flexibility to either improve the overall performance or increase the level of security by encrypting a greater number of code subsections while still meeting the original performance constraints.

*D. Advantages of our Approach*

Based on these two previous examples, we can summarize the advantages that our approach contains over current software protection methodologies as follows:

- 1) Our approach simultaneously addresses multiple attacks on software integrity by limiting both code understanding and code tampering.
- 2) The compiler’s knowledge of program structure, coupled with the programmability of the FPGA, provides the application developer with an extreme amount of flexibility in terms of the security of the system. The techniques available to any given application range from simple obfuscation that provides limited protection with a minimal impact on performance to a full cryptographic solution that provides the highest level of security with a more pronounced impact on performance.
- 3) Our approach provides the ability to combine both hardware specific techniques with hardware-optimized implementations of several of the software-based methods proposed recently [9], [8].
- 4) The selection of the FPGA as our secure component minimizes additional hardware design. Moreover, the choice of a combined processor/FPGA architecture enables our system to be immediately applicable to current SoC designs with processors and reconfigurable logic on-chip, such as the Xilinx Virtex-II Pro architecture [58].
- 5) As opposed to ASICs, FPGAs are especially well-suited towards cryptographic applications as their reconfigurable nature allows cipher designers to make changes to the implementation after the initial time of programming - this can be used to fix previously-unseen flaws in the cipher algorithm, switch between a set of parameters at run-time, or optimize the implementation for a predetermined fixed range of inputs.

We have up to this point in this paper demonstrated the usefulness of our architecture by providing specific examples of how our approach can be used in a software protection scheme. What remains to be seen, however, is to what extent the insertion of the FPGA in the instruction memory hierarchy

effects system security and performance. We address this question in the following section.

## V. ANALYZING PERFORMANCE OVERHEAD

Since the majority of the techniques we leverage operate on a single program basic block at a time, it makes sense to analyze the effect of the FPGA on instruction memory hierarchy performance at that level of granularity. We begin by considering the replacement penalty of a single block of cache directly from a pipelined main memory. With a fixed block size and memory bus width (in terms of number of bytes), we can estimate the penalty as the sum of an initial nonsequential memory access time for the first bytes from memory plus the delay for a constant amount of sequential accesses, which would be proportional to the product of the block size with the inverse of the memory bus width.

Now, considering a basic block of instructions of a certain length in isolation from its surrounding program, we can estimate the number of instruction cache misses as the number of bytes in the basic block divided by the number of bytes in the cache block, as each instruction in the basic block would be fetched sequentially. Consequently the total instruction cache miss delay for a single fetched basic block can be approximated as the product of the number of fetches and average fetch delay as described above.

What effect would our software protection architecture have on performance? We identify two dominant factors: the occasional insertion of new instructions into the executable and the placement of the FPGA into the instruction fetch stream. For the first factor, the inserted instructions will only add a small number of cache misses for the fetching of the modified basic block, since for most cases the number of inserted bytes will be considerably smaller than the size of the cache block itself. For the second factor, we note that the majority of the operations performed in the FPGA can be modeled as an increase in the instruction fetch latency. Assuming a pipelined implementation of whatever translation and validation is performed for a given configuration, we can estimate the delay for our FPGA as that of a similar bandwidth memory device, with a single nonsequential access latency followed by a number of sequential accesses. In the remainder of this section, we explain our experimental approach and then provide quantitative data that demonstrates how these terms are affected by the security requirements of any application.

### A. Experimental Methodology

For the following experiments, we incorporated a behavioral model of our software protection FPGA hardware into the SimpleScalar/ARM tool set [59], a series of architectural simulators for the ARM ISA. We examined the performance of our techniques for a memory hierarchy that contains separate 16 KB 32-way associative instruction and data caches, each with 32-byte lines. With no secondary level of cache, our non-sequential memory access latency is 100 cycles and our sequential (pipelined) latency is 2 cycles. We inserted our protective techniques into the back-end of a modified version of the `gcc` compiler targeting the ARM instruction set.

To evaluate our approach, we adapted six benchmarks from two different embedded benchmark suites. From the MediaBench [60] suite we selected two different voice compression programs: *adpcm* – which implements Adaptive Differential Pulse Code Modulation decompression and compression algorithms, and *g721* – which implements the more mathematically complex CCITT (International Telegraph and Telephone Consultative Committee) standard. We also selected from MediaBench the benchmark *pegwit*, a program that implements elliptic curve algorithms for public-key encryption. From the MiBench [61] embedded benchmark suite we selected three applications: *cjpeg* – a utility for converting images to JPEG format through lossy image compression, *djpeg* – which reverses the JPEG compression, and *dijkstra* – an implementation of the famous Dijkstra’s algorithm for calculating shortest paths between nodes, customized versions of which can be found in network devices like switches and routers.

It should be noted that both our simulation target and selected workload characterize a processor that could be found in a typical embedded system. This choice was motivated by the fact that commercial chips have relatively recently been developed that incorporate both embedded processors and FPGAs onto a single die (ex. Xilinx Virtex-II Pro Platform FPGAs [58]). Consequently, even though we see our techniques as one day being useful to a wide range of systems, the tradeoffs inherent in our approach can be best demonstrated through experiments targeting embedded systems, as these results are directly applicable to current technology.

### B. Initial Results

Using this simulation platform, we first explored the effects of our approach on performance and resultant security with an implementation of the tamper-resistant register encoding example from Section 3.2. In this configuration, the compiler manipulates sequences of register-based instructions to embed codes into the executable which are verified at run-time by the FPGA. As the operations required are relatively simple, for our experiments we assumed that the operations performed by the FPGA to decode individual instructions requires 1 clock cycle, and that verifying an individual basic block by comparing a function of the register sequence with a function of the filtered instruction requires 3 clock cycles.

As mentioned previously, it is often the case that the compiler will not have enough register-based instructions in a given basic block with which to encode a relatively long sequence string. Consequently, in these cases the compiler must insert some instructions into the basic block which encode a portion of the desired sequence but which otherwise do not affect processor state. However, as these “place-holder” instructions must also be fetched from the instruction cache and loaded in the processor pipeline they can cumulatively have a significant detrimental impact on performance.

While quantifying the security of any system is not a simple task, we can estimate the overall coverage of an approach independent of the instantiated tamper-checking computations. For our register encoding approach, we can measure this

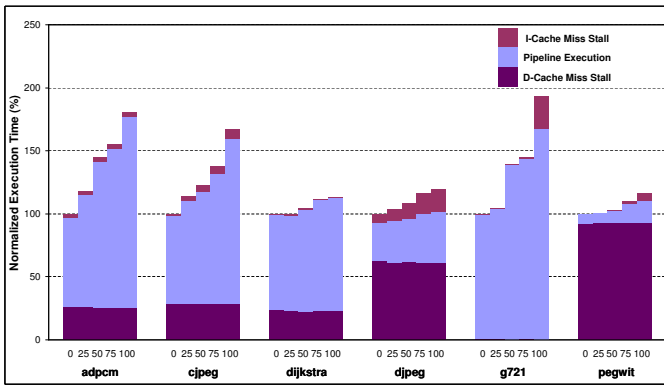


Fig. 10. Performance breakdown of the tamper-resistant register encoding scheme as a function of the basic block select rate

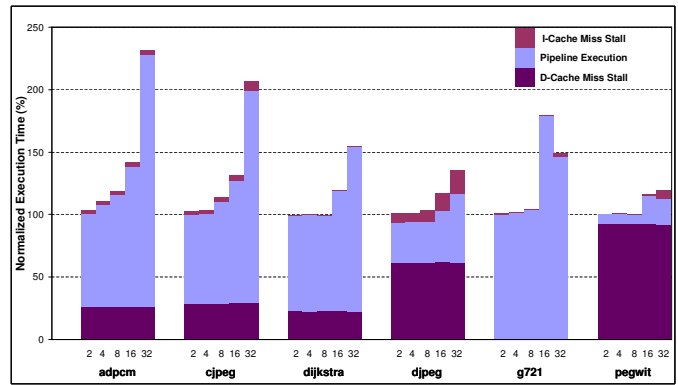


Fig. 11. Performance breakdown of the tamper-resistant register encoding scheme as a function of the register sequence length

aggressiveness as a function of both the desired register sequence length and the percentage of basic blocks that are protected. Fig. 10 considers the performance of our system when the sequence length is kept at a constant value 8 and the percentage of encoded basic blocks is varied from 0 - 100%. For each experiment the results are normalized to the unsecured case and are partitioned into three subsets: (1) the stall cycles spent in handling data cache misses, (2) the “busy” cycles where the pipeline is actively executing instructions, and (3) the stall cycles spent in handling instruction cache misses.

Several general points can be made about the results in Fig. 10. First, it is obvious that the selected embedded benchmarks do not stress even our relatively-small L1 instruction cache, as on average these miss cycles account for less than 3% of the total base case run-time. This is a common trait among embedded applications, as they can often be characterized as a series of deeply-nested loops that iterate over large data sets (such as a frame of video). This high inherent level of instruction locality means that, although the inserted register sequence instructions do affect the cycles spent in handling instruction cache misses, there is a relatively significant negative impact on the non-stall pipeline cycles. The average slowdown for this scheme is approximately 48% in the case where all the basic blocks are selected, but then drops to 20% if we only select half the blocks. These results clearly demonstrate the tradeoff between security and performance that can be managed in our approach.

In Fig. 11, we now consider the case where the basic block selection rate is kept constant at 25 and the effect of the desired register sequence length is examined. Although for most of the benchmark/sequence combinations the performance impact is below 50%, there are cases where lengthening the sequences can lead to huge increases in execution time, as the basic blocks are too short to supply the needed number of register-based instructions. This can be seen in the results for the most secure configuration (sequence length value of 32), where the average performance penalty is approximately 66%.

In our next set of experiments, we investigated the performance impact of the selective basic block encryption scheme. FPGA implementations of symmetric block often strive to optimize for either throughput or area. Recent implementations

of the Advanced Encryption Standard (AES) have reached a high throughput by unrolling and pipelining the algorithmic specification [52]. The resultant pipeline is usually quite deep. Consequently the initial access latencies are over 40 cycles, while the pipelined transformations can be executed in a single cycle. Assuming an AES implementation on our secure FPGA, we can make the intelligent assumption of a 50 cycle nonsequential access and a single cycle sequential access delay. Note that other AES implementations concentrate on optimizing other characteristics besides throughput (i.e. area, latency, throughput efficiency), but an examination of the performance impact of those designs on our architecture is outside the scope of this paper.

Fig. 12 shows the performance breakdown of our selective basic block encryption architecture as a function of the block select rate. The increased nonsequential access time for the FPGA has the expected effect on the instruction cache miss cycles, which for several of our benchmarks become a more significant portion of the overall execution profile. It is interesting to note that the average performance penalty of the case where the entire program is initially encrypted is less than 20%, a number that is significantly less than the seemingly simpler register-sequence encoding approach.

Why is this the case? This question can be answered by again examining the ratio of the instruction cache miss cycles to the pipeline execution cycles. Although the AES instantiation brings with it a significant increase in the former factor when compared to the register-based approach, the compiler pass that encrypts the executable only requires that two instructions (the start and stop message instructions) be inserted for every basic block. Consequently for applications that have excellent instruction cache locality such as our selected benchmarks, we would expect the performance of this scheme to be quite similar to that of the previous approach configured with a sequence length value of 2. A quick comparison of Fig. 10 with Fig. 12 shows this to be the case.

These results clearly demonstrate the flexibility of our approach. With the modification of a few compiler flags an application developer can evaluate both a tamper-resistant register encoding system that covers the entire application with a significant performance detriment, to a partially cov-

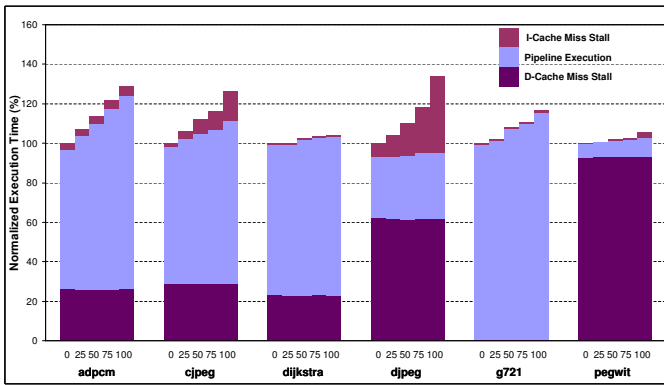


Fig. 12. Performance breakdown of the selective basic block encryption scheme as a function of the block select rate

ered cryptographic solution that has a much more measured impact on performance. Many of the more severe performance numbers are fully expected based on previous work - as an example, estimates for the XOM architecture [40] predict less than a 50% performance impact when coupled with a modern CPU. While the results in this section show that the instruction miss penalty isn't a dominant factor, this will generally not be the case when considering the larger applications that could be used with our approach. This, combined with the possibility that programmable FPGA resources will not all be allocated for software protection purposes motivates the configuration of these resources as performance-oriented architectural optimizations. We examine instruction prefetching in the following section to investigate the impact of such an approach.

## VI. ADDITIONAL PERFORMANCE OPTIMIZATIONS

For systems with a low instruction cache efficiency and long memory latencies, both buffering and prefetching have been shown to be useful optimizations for reducing cache misses and subsequently improving system performance. In this section we concentrate on prefetching and examine how including instruction prefetching logic on the FPGA affects performance.

Prefetching techniques tend to monitor cache access patterns in order to predict which instructions will be requested in the immediate future. These guesses are used as the basis whereby instructions are fetched from slow main memory into cache before they are requested. In general, the usefulness of a prefetching scheme is based on the percentage of cache misses that it is able to prevent. In other words, prefetching instructions will only improve performance if the average per-access overhead of the prefetching hardware is outweighed by the average overall fetch latency.

### A. FPGA-based Next-N-Line Prefetching

One of the most popular prefetching techniques is known as *next-N-line* prefetching. The main concept behind this technique is that when a block of instructions are requested main memory, the prefetching mechanism also returns the next N lines to the instruction cache. Our architecture for next-N-line prefetching can be described as follows (see Fig. 13).

A relatively small N-line buffer is instantiated on the FPGA which responds to instruction cache requests concurrently with the other security-related mechanisms of the FPGA. If the requested line is found, then the other FPGA operations are canceled and the data is returned directly from the buffer. If there is no match, then the result is returned from higher-level memory with no additional delay (assuming the latency in accessing higher-level memory is greater than that of the prefetching buffer). In either case, the next N lines following the requested address are then fetched through the remainder of the FPGA. This can be a quite simple operation for sequentially fetched blocks - it is very likely that many of the next N lines have already been recently prefetched and do not need to be fetched again. In order to minimize the average prefetching overhead we configured our FPGA to halt its current prefetching operation in the case of a new request from the L1 cache.

It is important to note that this requires no changes to the L1 instruction cache itself. This is nice feature as it means that if a cache miss leads to a poorly speculated prefetch operation (as could be the case in a procedure call that quickly returns to the callee location), the cache will not get polluted with useless lines. Also since our prefetch control is configured to not continue in its fetching of the next N lines in the case of a new request, such an architecture is likely to service any given request with a minimal additional delay.

### B. Results

We configured the FPGA in our selective basic block encryption approach to perform next-N-line instruction prefetching for the case when 25 percent of the eligible basic blocks are selected. The results for varying values of N can be seen in Fig. 14. Our most aggressive prefetching architecture (N = 16) does a fairly good job at predicting future misses - on average the prefetch buffer in this configuration successfully responds to 71% of the instruction cache requests. This correlates directly to an instruction cache miss cycle improvement. Fig. 14b shows that on average there is an average of 18% savings in instruction stall overhead. For our current benchmarks this translates to a relatively small overall performance improvement but given some benchmarks with a lower instruction locality this can have a great effect with no additional hardware cost.

## VII. CONCLUSIONS AND FUTURE WORK

Both the methods used by hackers in compromising a software system and the preventative techniques developed in response thereto have received an increased level of attention in recent years, as the economic impact of both piracy and malicious attacks continues to skyrocket. Many of the recent approaches in the field of software protection have been found wanting in that they either a) do not provide enough security, b) are prohibitively slow, or c) are not applicable to currently available technology. This paper describes a novel architecture for software protection that combines a standard processor with dynamic tamper-checking techniques implemented in

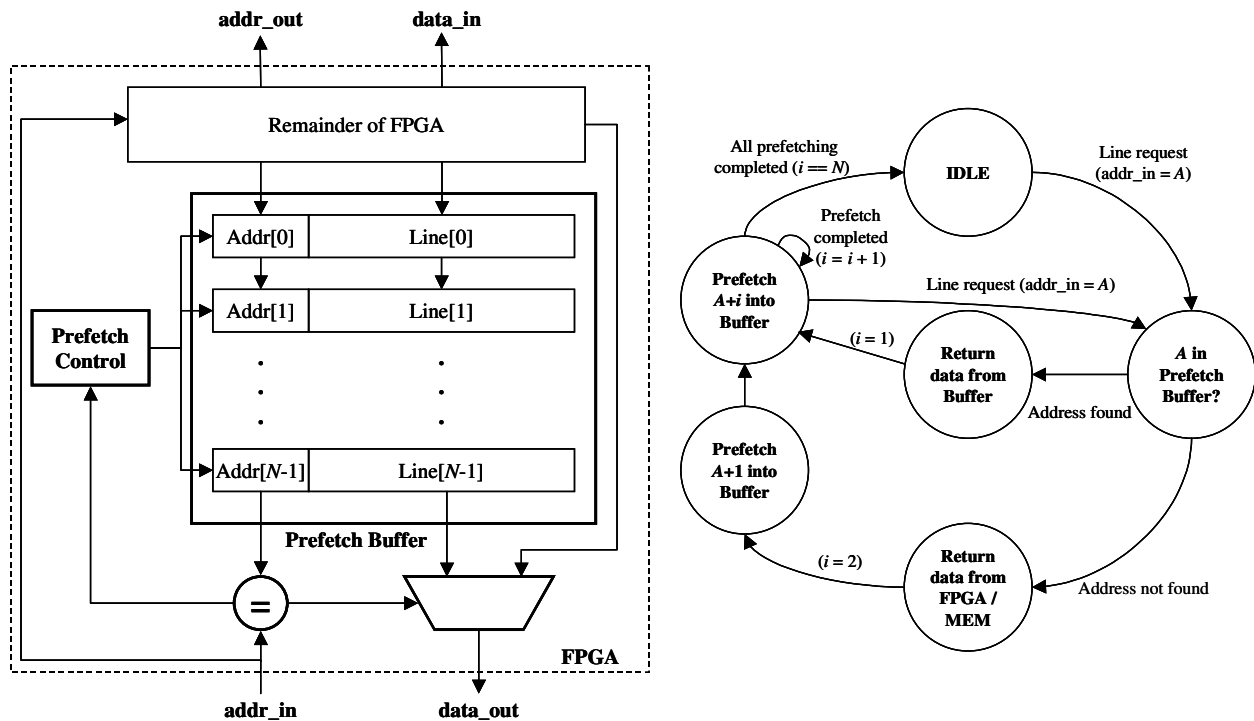


Fig. 13. Architecture and state diagram for a next-N-line prefetching technique implemented in FPGA

reconfigurable hardware. We have evaluated two distinct examples that demonstrate how such an approach provides much flexibility in managing the security/performance tradeoff on a per-application basis. Our results show that a reasonable level of security can be obtained through a combined obfuscating and tamper-proofing technique with less than a 20% performance degradation for most applications. When a highly-transparent solution is desired, FPGA resources not allocated for software protection can be used to mask some of the high latency operations associated with symmetric block ciphers.

Future work on this project will include implementing a method for calculating and storing hash values on the FPGA in order to secure data and file storage. Also, while our current simulation infrastructure is adequate for the experiments presented in this paper, it is also inherently limited in the sense that the component delays are completely user-defined. Although we intelligently estimate these values based on previous work in hardware design, in the future it would be considerably more convincing to assemble the results based on an actual synthesized output. This would require our compiler to be modified to generate designs in either a full hardware description language such as VHDL or Verilog, or at a minimum, a specification language that can drive pre-defined HDL modules. Taking this concept a step further, as our techniques can be fully implemented with commercial off-the-shelf components, it would be useful to port our architecture to an actual hardware board containing both a processor and FPGA. Such a system would allow us to obtain real-world performance results and also to refine our threat model based on the strengths and weaknesses of the hardware.

REFERENCES

- [1] International Planning and Research Corporation, "Seventh annual BSA global software piracy study," available at <http://www.bsa.org>, June 2002.
- [2] Computer Security Institute and Federal Bureau of Investigation, "CSI/FBI 2002 computer crime and security survey," available at <http://www.gocsi.com>, Apr. 2002.
- [3] C. Collberg and C. Thomborson, "Watermarking, tamper-proofing, obfuscation: tools for software protection," *IEEE Transactions on Software Engineering*, vol. 28, no. 8, pp. 735-746, Aug. 2002.
- [4] S. Chang, P. Litva, and A. Main, "Trusting DRM software," in *Proceedings of the W3C Workshop on Digital Rights Management for the Web*, Jan. 2001.
- [5] J. Wyant, "Establishing security requirements for more effective and scalable DRM solutions," in *Proceedings of the Workshop on Digital Rights Management for the Web*, Jan. 2001.
- [6] C. Collberg and C. Thomborson, "Software watermarking: models and dynamic embeddings," in *Proceedings of the 26th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, Jan. 1999, pp. 311-324.
- [7] R. Venkatesan, V. Vazirana, and S. Sinha, "A graph theoretic approach to software watermarking," in *Proceedings of the Fourth International Information Hiding Workshop*, Apr. 2001.
- [8] H. Chang and M. Atallah, "Protecting software code by guards," in *Proceedings of the ACM Workshop on Security and Privacy in Digital Rights Management*, Nov. 2000, pp. 160-175.
- [9] B. Horne, L. Matheson, C. Sheehan, and R. Tarjan, "Dynamic self-checking techniques for improved tamper resistance," in *ACM Workshop on Security and Privacy in Digital Rights Management*, Nov. 2001, pp. 141-159.
- [10] D. Aucsmith, "Tamper-resistant software: An implementation," in *Proceedings of the 1st International Workshop on Information Hiding*, May 1996, pp. 317-333.
- [11] A. W. Appel and E. W. Felten, "Proof-carrying authentication," in *Proceedings of the 6th ACM Conference on Computer and Communications Security*, Nov. 1999, pp. 52-62.
- [12] L. Bauer, M. A. Schneider, and E. W. Felten, "A proof-carrying authorization system," Department of Computer Science, Princeton University, Tech. Rep. CS-TR-638-01, Apr. 2001.
- [13] G. Necula, "Proof-carrying code," in *Proceedings of the 24th ACM*

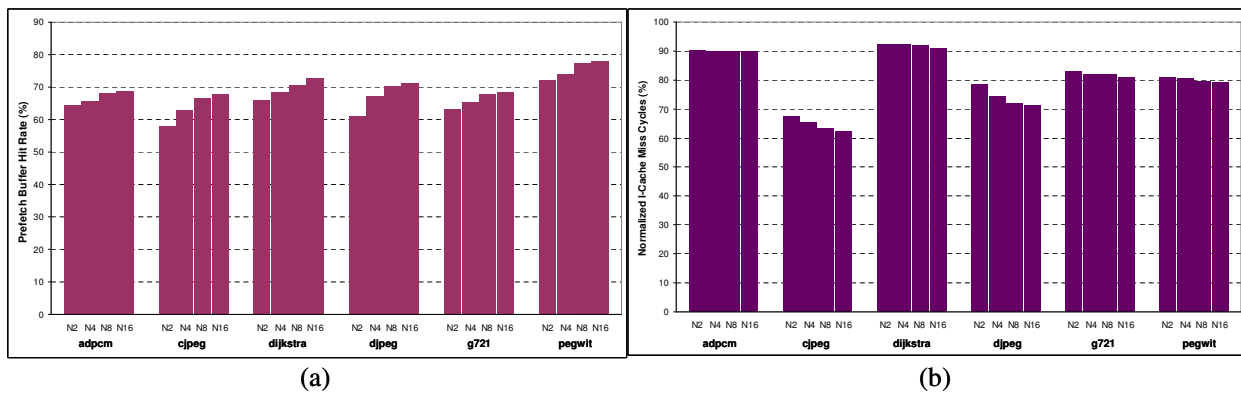


Fig. 14. Effectiveness of next-N-line instruction prefetching optimization as a function of N

SIGPLAN-SIGACT Symposium on Principles of Programming Languages, Jan. 1997, pp. 106–119.

[14] G. Necula and P. Lee, “Safe kernel extensions without run-time checking,” in *Proceedings of the 2nd USENIX Symposium on OS Design and Implementation*, Oct. 1996, pp. 229–243.

[15] D. Baifanz, D. Dean, and M. Spreitzer, “A security infrastructure for distributed Java applications,” in *Proceedings of the 2000 IEEE Symposium on Security and Privacy*, May 2000, pp. 15–26.

[16] C. Collberg, C. Thomborson, and D. Low, “A taxonomy of obfuscating transformations,” Department of Computer Science, The University of Auckland, Tech. Rep. 148, July 1997.

[17] —, “Breaking abstractions and unstructuring data structures,” in *Proceedings of the IEEE International Conference on Computer Languages*, May 1998, pp. 28–38.

[18] C. Wang, J. Davidson, J. Hill, and J. Knight, “Protection of software-based survivability mechanisms,” in *Proceedings of the 2001 IEEE/IFIP International Conference on Dependable Systems and Networks*, July 2001.

[19] C. Wang, J. Hill, J. Knight, and J. Davidson, “Software tamper resistance: obstructing the static analysis of programs,” Department of Computer Science, University of Virginia, Tech. Rep. CS-2000-12, May 2000.

[20] B. Barak, O. Goldreich, R. Impagliazzo, S. Rudich, A. Sahai, S. Vadhan, and K. Yang, “On the (im)possibility of obfuscating programs,” in *Proceedings of Advances in Cryptology (CRYPTO ’01)*, Aug. 2001, pp. 1–18.

[21] S. White, S. Weingart, W. Arnold, and E. Palmer, “Introduction to the Citadel architecture: Security in physically exposed environments,” IBM Research Division, T.J. Watson Research Center, Tech. Rep. RC 16682, May 1991.

[22] D. Tygar and B. Yee, “Dyad: A system for using physically secure coprocessors,” Department of Computer Science, Carnegie Mellon University, Tech. Rep. CMU-CS-91-140R, May 1991.

[23] B. Yee, “Using secure coprocessors,” Computer Science Department, Carnegie Mellon University, Tech. Rep. CMU-CS-94-149, May 1994.

[24] B. Yee and D. Tygar, “Secure coprocessors in electronic commerce applications,” in *Proceedings of the 1st USENIX Workshop on Electronic Commerce*, July 1995, pp. 155–170.

[25] S. Weingart, “Physical security for the mABYSS system,” in *Proceedings of the IEEE Symposium on Security and Privacy*, Apr. 1987, pp. 52–58.

[26] S. Weingart, S. White, W. Arnold, and G. Double, “An evaluation system for the physical security of computing systems,” in *Proceedings of the 6th Computer Security Applications Conference*, Dec. 1990, pp. 232–243.

[27] S. White and L. Comerford, “ABYSS: A trusted architecture for software protection,” in *Proceedings of the IEEE Symposium on Security and Privacy*, Apr. 1987, pp. 38–51.

[28] IBM, “Secure systems and smart cards,” available at [http://www.research.ibm.com/secure\\_systems](http://www.research.ibm.com/secure_systems), 2002.

[29] S. Smith, “Secure coprocessing applications and research issues,” Computer Research and Applications Group, Los Alamos National Laboratory, Tech. Rep. LA-UR-96-2805, Aug. 1996.

[30] S. Smith and V. Austel, “Trusting trusted hardware: towards a formal model of programmable secure coprocessors,” in *Proceedings of the 3rd USENIX Workshop on Electronic Commerce*, Aug. 1998, pp. 83–98.

[31] H. Gobiuff, S. Smith, D. Tygar, and B. Yee, “Smart cards in hostile environments,” in *Proceedings of the 2nd USENIX Workshop on Electronic Commerce*, Nov. 1996, pp. 23–28.

[32] B. Schneier and A. Shostack, “Breaking up is hard to do: modeling security threats for smart cards,” in *Proceedings of the USENIX Workshop on Smartcard Technology*, May 1999, pp. 175–185.

[33] P. Clark and L. Hoffman, “BITS: A smartcard protected operating system,” *Communications of the ACM*, vol. 37, no. 11, pp. 66–70, Nov. 1994.

[34] MIPS, “SmartMIPS application-specific extension,” available at <http://www.mips.com>, 2004.

[35] VIA, “Padlock hardware security suite,” available at <http://www.via.com>, 2004.

[36] Intel, “Intel LaGrande technology,” available at <http://www.intel.com>, 2004.

[37] ARM, “ARM TrustZone technology,” available at <http://www.arm.com>, 2004.

[38] Trusted Computing Group, <http://www.trustedcomputing.org>, 2003.

[39] Microsoft, “Next-generation secure computing base,” available at <http://www.microsoft.com>, 2004.

[40] D. Lie, C. Thekkath, M. Mitchell, P. Lincoln, D. Boneh, J. Mitchell, and M. Horowitz, “Architectural support for copy and tamper resistant software,” in *Proceedings of the 9th International Conference on Architectural Support for Programming Languages and Operating Systems*, Nov. 2000, pp. 168–177.

[41] D. Kirovski, M. Drinic, and M. Potkonjak, “Enabling trusted software integrity,” in *Proceedings of the 10th International Conference on Architectural Support for Programming Languages and Operating Systems*, Oct. 2002, pp. 108–120.

[42] X. Zhuang, T. Zhang, H.-H. Lee, and S. Pande, “Hardware assisted control flow obfuscation for embedded processors,” in *Proceedings of the International Conference on Compilers, Architectures and Synthesis for Embedded Systems (CASES)*, Oct. 2004.

[43] X. Zhuang, T. Zhang, and S. Pande, “HIDE: an infrastructure for efficiently protecting information leakage on the address bus,” in *Proceedings of the International Conference on Architectural Support for Programming Languages and Operating Systems*, Oct. 2004.

[44] A. Dandalis, V. Prasanna, and J. Rolim, “An adaptive cryptographic engine for IPsec architectures,” in *Proceedings of the IEEE Symposium on Field-Programmable Custom Computing Machines*, Apr. 2000, pp. 132–144.

[45] V. Prasanna and A. Dandalis, “FPGA-based cryptography for internet security,” in *Online Symposium for Electronic Engineers*, Nov. 2000.

[46] J.-P. Kaps and C. Paar, “Fast DES implementations for FPGAs and its application to a universal key-search machine,” in *Proceedings of the 5th Annual Workshop on Selected Areas in Cryptography*, Aug. 1998, pp. 234–247.

[47] R. Taylor and S. Goldstein, “A high-performance flexible architecture for cryptography,” in *Proceedings of the Workshop on Cryptographic Hardware and Software Systems*, Aug. 1999.

[48] Actel, “CoreDES data sheet, v2.0,” available at <http://www.actel.com>, 2003.

[49] —, “Design security with Actel FPGAs,” available at <http://www.actel.com>, 2003.

[50] O. Kommerling and M. Kuhn, “Design principles for tamper-resistant

- smartcard processors,” in *Proceedings of the USENIX Workshop on Smartcard Technology*, May 1999, pp. 9–20.
- [51] A. Elbirt, W. Yip, B. Chetwynd, and C. Paar, “An FPGA implementation and performance evaluation of the AES block cipher candidate algorithm finalists,” in *The Third Advanced Encryption Standard (AES3) Candidate Conference*, Apr. 2000, pp. 13–27.
- [52] K. U. Jarvinen, M. T. Tommiska, and J. O. Skytta, “A fully pipelined memoryless 17.8 Gbps AES-128 encryptor,” in *Proceedings of the International Symposium on Field Programmable Gate Arrays (FPGA)*, Feb. 2003, pp. 207–215.
- [53] J. Daeman and V. Rijmen, “The block cipher Rijndael,” in *Smart Card Research and Applications*, ser. Lecture Notes in Computer Science, J.-J. Quisquater and B. Schneier, Eds. Springer-Verlag, 2000, vol. 1820, pp. 288–296.
- [54] C.-K. Luk and T. Mowry, “Cooperative prefetching: Compiler and hardware support for effective instruction prefetching in modern processors,” in *Proceedings of the 31st Annual International Symposium on Microarchitecture (MICRO-31)*, Dec. 1998.
- [55] G. Reinman, B. Calder, and T. Austin, “Fetch directed instruction prefetching,” in *Proceedings of the 32nd Annual International Symposium on Microarchitecture (MICRO-32)*, Nov. 1999.
- [56] I.-C. K. Chen, C.-C. Lee, and T. Mudge, “Instruction prefetching using branch prediction information,” in *Proceedings of the International Conference on Computer Design (ICCD)*, Oct. 1997, pp. 593–601.
- [57] J. Zambreno, A. Choudhary, R. Simha, and B. Narahari, “Flexible software protection using HW/SW codesign techniques,” in *Proceedings of Design, Automation, and Test in Europe*, Feb. 2004, pp. 636–641.
- [58] Xilinx, “Virtex-II Pro Platform FPGA data sheet,” available at <http://www.xilinx.com>, 2003.
- [59] D. Burger and T. M. Austin, “The simplescalar tool set, version 2.0,” Department of Computer Science, University of Wisconsin-Madison, Tech. Rep. CS-TR-97-1342, June 1997.
- [60] C. Lee, M. Potkonjak, and W. H. Mangione-Smith, “MediaBench: A tool for evaluating and synthesizing multimedia and communications systems,” in *Proceedings of 30th Annual International Symposium on Microarchitecture*, Dec. 1997, pp. 330–335.
- [61] M. R. Guthaus, J. S. Ringenberg, D. Ernst, T. M. Austin, T. Mudge, and R. B. Brown, “MiBench: A free, commercially representative embedded benchmark suite,” in *Proceedings of the 4th IEEE Annual Workshop on Workload Characterization*, Dec. 2001, pp. 3–14.