fASLR: Function-Based ASLR via TrustZone-M and MPU for Resource-Constrained IoT Systems

Lan Luo¹⁰, Xinhui Shao, Zhen Ling⁹²T 1.4 2.2169 1.75.1 Member, IEEE

Yumeng Wei, and Xinwen Fu, Senior Member, IEEE

Abstract—The address space layout randomization (ASLR) has been widely deployed on modern operating systems against code reuse attacks (CRAs), such as return-oriented programming (ROP) and jump-oriented programming (JOP). However, porting ASLR to resource-constrained IoT devices is a great challenge due to the limited memory space for randomization. We propose a function-based ASLR scheme (fASLR) for IoT runtime security utilizing the ARM TrustZone-M technology and the memory protection unit (MPU) supported by ARM Cortex-M processors. fASLR loads a function from the flash and randomizes its base address in a randomization region in RAM when the function is being called. We design novel mechanisms on cleaning up finished functions from the RAM and memory addressing to tackle the complexity of function relocation and randomization. Optimizations are applied to effectively reduce overhead introduced by runtime memory management. We also formally prove that user applications will run correctly with fASLR enabled. Compared with the related work, a prominent advantage of fASLR is that fASLR can run an application even if the application code cannot be completely loaded into RAM for execution. We test fASLR with 21 applications. The experimental results show that fASLR achieves a high randomization entropy and incurs a runtime overhead of less than 10%.

Index Terms—Address space layout randomization (ASLR), code reuse attacks (CRAs), Internet of Things, microcontroller, TrustZone.

Manuscript received 30 May 2022; revised 18 June 2022; accepted 26 June 2022. Date of publication 13 July 2022; date of current version 7 September 2022. This work was supported in part by the National Key Research and Development Program of China under Grant 2018YFB2100300; in part by the National Natural Science Foundation of China under Grant 62022024, Grant 61972088, Grant 62072103, Grant 62102084, Grant 62072102, Grant 62072098, and Grant 61972083; in part by the U.S. National Science Foundation (NSF) under Award 1931871 and Award 1915780; in part by the U.S. Department of Energy (DOE) under Award DE-EE0009152; in part by the Jiangsu Provincial Natural Science Foundation for Excellent Young Scholars under Grant BK20190060; in part by the Jiangsu Provincial Natural Science Foundation of China under Grant BK20190340; in part by the Jiangsu Provincial Natural Science foundation Security under Grant BM2003201; in part by the Key Laboratory of Computer Network and Information Integration of Ministry of Education of China under Grant 93K-9; and in part by the Collaborative Innovation Center of Novel Software Technology and Industrialization. (*Corresponding author: Zhen Ling.*)

Lan Luo is with the Department of Computer Science, University of Central Florida, Orlando, FL 32816 USA (e-mail: lukachan@knights.ucf.edu).

Xinhui Shao and Yumeng Wei are with the School of Cyber Science and Engineering, Southeast University, Nanjing 210096, China (e-mail: xinhuishao@seu.edu.cn; yumeng5@seu.edu.cn).

Zhen Ling and Huaiyu Yan are with the School of Computer Science and Engineering, Southeast University, Nanjing 210096, China (e-mail: zhenling@seu.edu.cn; huaiyu_yan@seu.edu.cn).

Xinwen Fu is with the Department of Computer Science, University of Massachusetts Lowell, Lowell, MA 01854 USA (e-mail: xinwenfu@cs.uml.edu).

Digital Object Identifier 10.1109/JIOT.2022.3190374

I. INTRODUCTION

Windows, macOS, Lniux, Android, and iOS.

In this article, we focus on defending against CRAs for resource-constrained IoT devices, particularly those running on microcontrollers (MCUs). It is an intuitive idea to port existing security schemes to IoT platforms. We study the use of address space layout randomization (ASLR) in memoryconstrained IoT devices to mitigate CRAs, such as the returnoriented programming (ROP) and jump-oriented programming (JOP) by randomizing the memory layout of code and data. Modern operating systems often implement the following ASLR scheme. When an executable is loaded into RAM, its base (start) address is randomly chosen while the executable structure is kept almost intact. Fine-grained ASLR strategies have been proposed and randomize executable code at fine levels of basic blocks, functions, or instructions [3] within a loaded application image. However, porting ASLR to resourceconstrained IoT devices is a great challenge due to the limited memory space.

We propose a function-based ASLR scheme (fASLR) based on the ARM Cortex-M processor with TrustZone-M enabled [4] to protect MCU-based IoT devices from CRAs that require the knowledge of locations of executable code snippets, such as ROP and JOP. fASLR takes advantage of hardwarebased isolation provided by TrustZone-M. The runtime fASLR is located in a trusted execution environment (TEE), namely, the secure world (SW) of TrustZone, while the application code is in a rich execution environment (REE), namely, the nonsecure world (NSW), and denoted as the NS app. The NS app is protected by the memory protection unit (MPU). When a function in the MPU protected region is called, a hardware exception is raised and the control flow of the function call is redirected to our custom exception handler, which is used by the runtime fASLR for callee randomization. Compared to the

(327-4662 © 2022 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See https://www.ieee.org/publications/rights/index.html for more information. most recent related work [5] that requires loading the whole application code into RAM, fASLR can run an application even if the application on flash is too large to be completely loaded into RAM.

fASLR is user friendly and does not require any code instrumentation for the user code. A programmer only needs to compile the NS code via the GCC compiler with specific compiler flags. fASLR implements a block-based memory management strategy to manage randomized functions in RAM. To reduce overheads introduced by runtime fASLR, three optimizations are adopted: 1) fASLR cleans up finished functions from RAM only when the randomization region (RR) is full; 2) a novel stack unwinding mechanism is devised to precisely find all functions that are safe to be cleaned; and 3) function call rewriting is used to replace the destination addresses of call instructions with the base addresses of the corresponding randomized loaded callees so that the call instructions can jump directly to the target loaded callees without raising exceptions.

A conference version of this article [6] mainly focuses on the optimized fASLR. In this journal version, we add a basic memory management strategy and compare it with the optimized fASLR in the conference paper. In addition, we formally prove mechanisms adopted by fASLR do not affect execution correctness of the NS app. We also port three test apps with relatively high time overheads evaluated in the conference version to a more powerful TrustZone-M-enabled MCU and compare the overheads at different MCUs. Finally, we discuss the compatibility of using fASLR with other protection mechanisms for securing the runtime execution of MCU-based IoT devices.

Our major contributions are summarized as follows.

- We propose a *function-based ASLR* scheme for resourceconstrained IoT devices with limited RAM and flash. fASLR dynamically loads only needed functions into RAM and randomizes their entry addresses so as to achieve large randomization entropy.
- 2) Novel schemes are designed for fASLR to perform memory management and addressing. We carefully address the issue of addressing since functions are randomly moved around. Finished functions are removed from RAM when there is no RAM to execute new function calls. Therefore, our scheme can run an NS app that is larger than the RAM.
- 3) We formally prove that the NS app still runs correctly with fASLR via logical reasoning.
- 4) We implement fASLR with a TrustZone-M-enabled MCUs, SAM L11, and STM32. We validate the feasibility and performance of fASLR with 21 applications on SAML11 and three applications on STM32. fASLR incurs a runtime overhead of less than 10% for all the applications. We also compare the performance of fASLR on SAM L11 and STM32 and show larger RAM can reduce the overhead as expected.

Roadmap: The remainder of this article is structured as follows: we first discuss the background of TrustZone-M-enabled processors in Section II. The threat model, design goals, and system architecture of ASLR are then presented in Section III. We also demonstrate the workflow of fASLR and two technical challenges in this section. In Section IV, we discuss the technical challenges and present our solutions. We prove the execution validity of the NS app with fASLR in Section V. In addition, we analyze the effectiveness and performance of fASLR in Section VI, and present experimental results in Section VII. Finally, we discuss the compatibility of fASLR with other security mechanisms in Section VIII, present related work on fine-grained ASLR techniques for embedded systems in Section IX, and conclude this article in Section X.

II. BACKGROUND

In this section, we introduce ARM Cortex-M MCUs and TrustZone-M, which is used in this article. ARM Cortex-M is a series of processors optimized for MCUs. Such processors come equipped with MPU, specific exception model, and different processor modes for security concerns.

Memory Protection Unit: The MPU is the security extension that enforces memory access permissions (i.e., read, write, and execute) for memory regions. Any access violation at memory address protected by MPU will trigger the ARM HardFault exception handling. Once such an exception is triggered, the processor will stop the current execution and execute the exception handler in the SW to respond to the exception. Before the execution of the exception handler, the processor context is first preserved in the call stack. The stack frame of the exception context is composed of the status registers (xPSR), program counter (PC), link register (LR), and general-purpose registers R12 and R0 to R3. At the same time, LR is set with EXC_RETURN, which is the address where the exception occurs.

Processor Mode: ARM Cortex-M processors support two processor modes, i.e., thread mode and handler mode. While the thread mode is for normal program execution and can be either privileged or unprivileged, the handler mode is for exception handling and only supports privileged software execution.

TrustZone-M-Enabled MCU: MCU often runs either a baremetal-embedded application that usually consists of an infinite loop performing a sequence of operations or a lightweight real-time operating system (RTOS) such as FreeRTOS [7]. The applications or RTOS are stored in the MCU on-chip memory. There are two types of MCU on-chip memory: 1) flash or EEPROM as the nonvolatile memory and 2) RAM as the volatile memory. Usually, an MCU program is programmed into the flash and executed directly in the flash, though MCUs allow running code snippets in RAM for performance concerns.

TrustZone-M is a hardware-based security technique designed for MCUs, providing two isolated execution environments named SW as the TEE, and NSW as the REE. The on-chip resources, such as memories and peripherals are divided into the two worlds as well. For simplicity, we use the word "Secure" to describe resources in the SW and use "Nonsecure" for those belonging to the NSW. Secure application (abbreviated as app), for example, is an app in the SW. In a TrustZone-enabled system, an SW program is able to access resources in both worlds, while a program in the NSW is considered to be untrusted and can only access resources in the NSW directly.

III. fASLR-ENABLED SYSTEM

In this section, we first present the threat model and design goals of our ASLR scheme—fASLR. We then introduce the architecture of an fASLR-enabled system and the workflow of fASLR. Finally, we discuss challenges for implementing a practical fASLR.

A. Threat Model

fASLR leverages ARM Cortex-M processors and hardwarebased techniques, including TrustZone-M, MPU, and exception handling mechanism. Based on the hardware isolation provided by TrustZone-M, on-device system resources are divided into two worlds, namely, the SW and the NSW.

We assume a TrustZone-M-enabled device has the following security features.

- Main components of fASLR reside in the SW and can be fully trusted. The application (denoted as NS app) is located in the NSW and may be vulnerable.
- The NS app is located at a fixed address in the NS flash and is executed from the flash (instead of RAM) by default.
- 3) The device supports the memory protection mechanisms such as the MPU.
- We assume an adversary has the following capabilities.
- The NS app may be subject to CRAs such as the ROP attack.
- The adversary can obtain the binary of the NS app, disassemble the binary, and obtain code gadgets for CRAs.

B. Design Goals

fASLR is designed to achieve the following goals.

- 1) *Mitigating CRAs:* The scheme shall provide dynamic function-level code randomization for resource-constraint IoT devices to mitigate CRAs, which require a certain chain of gadgets found in the NS app. The randomization shall achieve high entropy to defeat brute-force guessing attacks.
- 2) *Usability:* The scheme shall be user friendly and will not add much burden of programming.
- 3) *Low Runtime Overhead:* The proposed scheme cannot introduce large overhead in terms of time and space and affect the NS app performance much.

C. System Architecture

As illustrated in Fig. 1, fASLR has three key components: 1) the *static preprocessing module* (SPM) for compilation time preparation; 2) the *boot engine* (BE) for boot time configuration; and 3) the *function randomization engine* (FRE) for runtime function-level randomization.

Static Preprocessing Module: The SPM serves two major purposes.





	Pointer (ptr) Size	Payload Padding bytes		
1	ptr	Unused space	ptr Function 1 ptr	Unused

Function 2

Function 1

ptr

of fALSR with respect to the timing of function cleaning and calling loaded functions.

- 1) *Call Stack Unwinding:* Finished functions are found through unwinding the Nonsecure call stack.
- Cleaning on Demand: Finished functions are cleaned up only if the available RR space is not large enough for the callee.
- 3) *Call Instruction Rewriting:* We further reduce the runtime overhead by overwriting a call instruction in a loaded function if the callee of that call instruction has already been loaded into RAM.

Call Stack Unwinding: The key of function cleaning is to distinguish finished functions from all loaded functions in RAM. However, it is difficult to trace all finished functions at runtime because fASLR runtime does not capture any function return information. Instead, our approach finds ancestor functions of the current callee, and records all loaded functions. Any function that is a loaded function but not an ancestor function is a finished function that can be disposed. Now, the problem is decomposed to record all loaded functions and find all ancestor functions.

Like the trace stack working for the baseline memory man-

TABLE I ATOMIC OPERATIONS OF RUNTIME FASLR

	Momie of Ekanons of Konfine MoEk	
Operation A: Act-1 \rightarrow Act-2A \rightarrow	3]
Operation B: Act-1 \rightarrow Act-2B \rightarrow Operation C: Act-1 \rightarrow Act-2C \rightarrow	, 3	
Action-1 Obtaining the base a	of the callee in flash and querying the corresponding function record from the Loading Queue	
Action-2A Condition If the fun		



ongoing fASLR operation with regard to the trapped function call. A future function call is one that will occur after the active function call. A function return will definitely occur after the active function call is executed.

As an fASLR operation performs when an active function call occurs, the operation is in full control of the execution of the active function call. The validity of the active function call depends on how the operation affects the program logic of the call. On the other side, any future function call that may occur after the current operation is generated by the execution of a call instruction. The validity of such a call depends on whether the call instruction points to the correct callee function. In other words, the memory contents at the location pointed by the destination address of the call instruction (i.e., static call information) must be the correct callee function (either the original callee in flash or the corresponding callee duplicate in RR).

A function return obtains the return location through return address. In ARM, leaf subroutine and nonleaf subroutine use different ways to obtain the return address. For a leaf subroutine, the return address is stored in the lr register, while a nonleaf subroutine uses the return address stored in the stack frame. Because fASLR does not intervene in the process of putting a return address into the lr register or onto a stack frame, it can be affirmed that all return addresses are generated correctly as the original NS app (without fASLR) does. After the generation of the return addresses, fASLR never changes the lr register and stack frames used by normal program execution. So, until a function returns, the return address, no matter it is in 1r or on the stack, remains unchanged. The validity of future function returns therefore requires that the memory content pointed by the return addresses is valid for function returns.

In summary, the validity of the active function call depends on if the program logic is affected by an fASLR operation; the validity of future calls is determined by static call information and content validity of RR; Similarly, the validity of future returns relies on the content validity of RR. In addition, the auxiliary data structures, as we introduced in concepts, must remain valid throughout the program execution so that any fASLR operation can perform as designed. Thus, we can conclude that the validity of function calls and returns critically relies on the validity of the following four objects throughout the program execution.

- 1) Program logic of the active call.
- 2) Static call information.
- 3) Auxiliary data structures (LQ and RL).
- 4) Randomization region.

We name the union of these four objects as the *critical reliance set*. The validity of the critical reliance set is in fact the necessary and sufficient condition for the validity of the function calls and returns. In later proof, we check the validity of the critical reliance set whenever the validity of function calls and returns need to be verified.

3) Assumptions:

Assumption 1 (Correct Initial System Status): When the system starts, the initial status of the whole system, including

the program logic, static call information, auxiliary data structures, and memory state, is correct.

Assumption 2 (Serial Execution System): The MCU that runs the program with fASLR is a serial execution system, in which one computation can begin only after the previous computation completes without parallelization.

D. Propositions

Proposition 1: Any operation of fASLR affects only function calls and returns of the NS app execution.

Proof: Program logic, which refers to the implementation of the program's design, is mutually determined by control flow and data flow at runtime. For the operations listed in Table I, the runtime fASLR does not modify any data flow during execution. We thus focus on the influence on the control flow. The control flow of a program can be divided into control flow within a function (i.e., branches inside a function and nonbranch execution) and control flow between functions (i.e., function calls and returns). It can be seen that operations listed in Table I do not affect any execution within a function. Hence, all Operations A, B, and C may affect only function calls and returns of the NS app execution.

E. Lemmas and Theorem

Lemma 1: The correctness of the NS app execution maintains when Operation A finishes.

Proof: According to Proposition 1, Operation A can only affect function calls and returns of the NS app. The validity of function calls and returns, as we introduced in concepts, can be verified by checking the validity of the critical reliance set. Therefore, we prove Lemma 1 through justifying that the critical reliance set remains valid when Operation A finishes. Since Operation A is the ordered combination of three actions, we first analyze whether each action affects the validity of the critical reliance set is valid upon the entry of Operation A. Such an assumption is natural and common, and will be consistently used among the proofs of all lemmas.

Action-1 only reads the entry address of the callee from the stack and search the function record from the LQ. It never writes any values or memory contents, and would not affect the validity of the critical reliance set.

The condition in Action-2A first guarantees the callee can be loaded into the RR. The loading action changes the status of the RR and LQ. Note that before this action, both RR and LQ are valid. So, we focus on the changes applied on them. After loading the callee (denoted as x) into the RR, for the occupied memory region (denoted as m) of the callee, indeed a new function record {Addr_x, Addr'_x, S_x } is created in the LQ to record that this region is being used. These are the solely changes to the RR and LQ, and these two changes are entirely correspondent. Conditions for the validity of RR are satisfied. Therefore, both LQ and RR remain valid.

When Action-3 is applied, the RR contains x, which is the duplicate of the callee. Action-3 forwards the control flow to this duplicate so the program logic is exactly the same as before.

So far, we have proved that any action of *Operation A* does not affect the validity of the critical reliance set. Hence, the critical reliance set will remain valid at the exit of Operation A. The correctness of the NS app execution holds at the exit of Operation A.

Lemma 2: The correctness of the NS app execution maintains when Operation B finishes.

Proof: Compared to Operation A, the only difference of *Operation B* is that it performs Action-2B instead of Action-2A. Therefore, we focus on the changes brought by Action-2B. The effects of other actions of Operation B are the same as Operation A.

Action-2B solely changes the static call information, i.e., the destination address of the active function call and adds the corresponding rewriting record to the RL. The new destination address $Addr'_y$ obtained from LQ is the correct base address of the duplicate callee because LQ is valid upon the entry of this operation. So, it is straightforward to see that both the static call information and RL remain valid. The correctness of the NS app execution is kept when Operation B finishes.

Lemma 3: The correctness of the NS app execution holds when Operation C finishes.

Proof: We focus on Action-2C since it is the only difference between Operations C and A.

Action-2C involves function cleaning and call instruction restoring. Function cleaning changes RR and LQ. As for the memory content of RR, Action-2C solely cleans the memory with state O from the status view according to LQ and stack, and deletes corresponding function records in LQ. This means LQ remains valid, and the state transition of such memory region is $O \rightarrow A$, which follows the previously defined transition rule. The consistency between LQ and RR is kept and satisfies the first requirement for the validity of RR. Based on the function cleaning process, the status view of the RR before and after function cleaning can be presented as $\{I|O|A\}$ and $\{I'|O'|A'\}$, where I' = I, O' = Empty, and A' = A + O. Similarly, from the usage view, we have N' = N + O, U' = U - O. Recall the second condition of a valid RR ensures N = A and U = I + O. After cleaning, we have

$$U' = U - O = I + O - O = I = I' + O'$$

 $N' = N + O = A + O = A'.$

This means that the RR after cleaning satisfies the second condition as well. Now, we can conclude that RR and LQ are still valid after cleaning.

Function cleaning may affect static information and RL. When loaded functions are cleaned from RAM, rewritten call instructions with the destination addresses pointing to the cleaned functions need to be restored to point to their original callees in flash. This is exactly what we do in this action. While cleaning a function f, by scanning all corresponding function records in RL, fASLR can precisely identify in static information of the set of call instructions pointing to the cleaned function. For each found function record {Addr_f : Addr_i, Addr_f}, fASLR deletes the rewriting record and restores the instruction *i* to use the original address Addr_f of the callee as the destination address. So, RL remains



Fig. 7. Execution model of the NS app with fASLR.

accurate and all static call information remains valid after Action-2C.

We can conclude that the critical reliance set affected by Operation C remains valid. So, the NS app can execute correctly when Operation C is applied.

Lemma 4: The correctness of the NS app execution after any operation in runtime fASLR will hold until the next occurrence of a runtime fASLR operation, regardless of the in-between program execution.

Proof: So far, we have proved that the occurrence of any operation in runtime fASLR does not affect the validity of the critical reliance set. Hence, the correctness of the program execution holds at the exit of each runtime fASLR operation. Because only an operation in fASLR could possibly change the validity of the critical reliance set, such validity after an operation in fASLR holds until the next operation happens. Such validity will not be affected by the specific program execution during these two operations either because the critical reliance set cannot be changed by any program execution.

Recall that the validity of the critical reliance set is equivalent to the validity of the function calls and returns and, thus, equivalent to the overall correctness of the program execution. Based on the observations above and Lemmas 1–3, the validity of the critical reliance set prevails between consecutive runtime fASLR operations. We can deduce that the correctness of the NS app execution, not only holds after any operation in runtime fASLR but also holds until the right next occurrence of a runtime fASLR operation. Lemma 4 is proved.

Theorem 1: fASLR does not affect the correctness of the NS app execution.

Proof: As we analyzed at the beginning of this section, the execution of the NS app with runtime fASLR can be seen as inserting several fASLR operations into the original execution of the NS app, as it runs in a serial execution system as we assume in Assumption 2. Such an execution model is illustrated in Fig. 7. For each operation, we combine the operation with the normal program execution right after it, until the occurrence of the next operation, as an execution block. So, the whole program execution can be seen as a chain of such execution blocks.

We have proved that the NS app execution will remain correct when Operation A, B, or C completes in Lemmas 1–3 separately, and proved that the correctness of the NS app execution after any operation prevails until the right next operation occurs in Lemma 4. According to Assumption 1, the whole system is initialized correctly. Thus, the NS app execution is correct at the beginning of block 1 in Fig. 7, and remains correct at the end of this block, which is also the beginning of block 2. Similarly, the correctness of the program execution prevails until block 2 completes. The correctness will not be affected by specific operations or specific normal executions. Therefore, no matter what operations and what normal executions compose the execution blocks after block 2, the program execution will remain correct until the program execution finishes. In other words, fASLR does not affect the correctness of the NS app execution. Theorem 1 is proved.

VI. SECURITY AND PERFORMANCE ANALYSIS

In this section, we first analyze the effectiveness of fASLR against ROP, a representative CRA. Entropy is computed to quantify the randomness of gadgets required for the ROP attack, which indicates the difficulty of guessing the gadget locations in a brute-force way. We also study time and memory overheads introduced by fASLR.

A. Effectiveness Against ROP

The prerequisite of ROP is that the adversary knows where the ROP gadgets are. In an fASLR enabled system, an adversary can only use ROP gadgets in randomized functions relocated to the RR. Gadgets in the NS app stored in flash are nonexecutable, so it is hard for adversaries to use them. Recall that any MPU violation triggers the HardFault exception. As discussed in Section III-C, the FRE validates the return address of the exception by using the FT. Therefore, the FRE is incapable of identifying exceptions triggered by a ROP attack if the adversary targets the entry point of a function since normal function calls will trigger such exceptions as well. In other words, the adversary will succeed in reusing a whole function as a gadget for ROP attack. However, such gadgets are often of very low quality [11], [12] containing too many instructions. It is almost impossible for an adversary to assemble a chain of gadgets with such low-quality gadgets to achieve certain malicious goal.

An adversary may also guess the addresses of randomized functions in a brute-force way. However, our runtime randomization approach rebases a function every time as long as it has not been loaded into RAM and achieves high randomization entropy as analyzed below.

B. Randomization Entropy

fASLR mitigates the brute-force guessing attack as follows.

- fASLR restricts the number of functions that can be reused at a time. This is achieved by configuring the whole app image as nonexecutable. The only code snippets that can be utilized are functions relocated to the RR in the RAM.
- 2) Even if all the required gadgets can be found from the relocated functions, the adversary has to guess locations of all those functions at once. Formula (1) gives the total number (denoted as *C*) of possible function layouts in the RR

$$C = k! \binom{V+k}{k} \tag{1}$$

$\left f_2 \right \left f_3 \right $	
--	--

rebase the callee; 4) function loading, which reads and writes the function body; and 5) function rewriting, which overwrites the destination of the call instruction with the entry point of loaded function.

D. Memory Overhead

The components of fASLR deployed in the SW include the BE code, FRE code, FT, LQ, and RL. The FT is a static table with three 4-byte attributes and its size is linear to the total number of functions in the NS app. The LQ and RL are dynamic data structures that contain function records and rewriting records, respectively. Each function record has four 4-bytes and one 1-byte metadata, and a rewriting record contains four 4-bytes data. The maximum number of records that the LQ may use at runtime is equal to the number of functions in the NS app, while the maximum number of rewriting records in the RL is the total number of call instructions. Formula (4) presents the size of the FT (i.e., MO_t), LQ (i.e., MO_q), and RL (i.e., MO_l)

$$MO_t = N_f \times 3 \times 4 = 12N_f \tag{4}$$

$$MO_q = N_f \times (4 \times 4 + 1) = 17N_f$$
 (5)

$$MO_l = N_c \times 4 \times 4 = 16N_c \tag{6}$$

where N_f is the number of functions in the NS app, and N_c is the number of function calls in the NS app.

E. Size Requirement of the Randomization Region

fASLR will run out of memory (OOM) if a new function cannot fit into the RR and no function can be trimmed. To avoid such an OOM issue, there is a size requirement of the RR for a certain application. We define call path size as the total size of all functions on a call path. The RR should be no less than the largest call path of the application when fragmentation compaction is applied by the memory management scheme. We can calculate the size requirement by statically analyzing the application code and perform defragmentation to the RR if needed.

VII. EVALUATION

In this section, we first present the experimental setup. We then present the evaluation of randomization entropy, runtime overhead, and memory overhead.

A. Experiment Setup

fASLR is implemented and deployed on the SAM L11 Xplained Pro Evaluation Kit, a MCU development board using the ARM Cortex-M23 core with TrustZone-M enabled. SAM L11 has a 64-kB flash and a 16-kB SRAM.

Software in SAM L11 is built with the GNU Arm Embedded Toolchain. User code, namely, the NS app code, is compiled with two flags, *-mlong-calls* and *-fno-jump-tables*, to eliminate instructions using relative addressing. We recompile the C library with the same compiler flags. A Python script runs during the compilation time to collect function metadata and saves them in the FT. fASLR program and the FT are



Fig. 9. Air quality monitoring device.



Fig. 10. Entropy distribution.

part of the Secure application placed in the SW flash, while the user app is deployed in the NSW flash.

We evaluate the performance of fASLR with 21 applications, including our own air quality monitoring system (*AirQualityMonitor*). The air quality monitoring device, as shown in Fig. 9, consists of a SAM L11 development board, a PMSA003 air quality sensor module, and a SIM7000 cellular module. The NS app in SAM L11 periodically receives air quality data from PMSA003 and sends the data to SIM7000, which then transfers the data to the AWS IoT platform via secure MQTT protocol. The other 20 apps, including the CoreMark benchmark [13], two microbenchmarks *Cache Test* and *Matrix Multiply* created based on [14], nine benchmarks of BEEBS (with the prefix *Beebs*-) [15], and eight SAM L11 demo apps (with the prefix *AS*-) obtained from Atmel Start [16].

B. Randomization Entropy

The entropy of function randomization changes dynamically when a function call occurs. We explore the entropy for all test applications. For each measured pair of k and V, we calculate the corresponding entropy of function randomization according to (1) and (2). Fig. 10 is the box plot demonstrating the entropy distribution for each app. The smallest average entropy is around 80 which is still considered to be large enough to defend against brute-force guessing.

C. Runtime Overhead

fASLR introduces runtime overhead since it intercepts every function call of the NS app for function randomization. We evaluate the time overhead by measuring and comparing the execution time of an application with and without fASLR. We use the internal systick timer of the Cortex-M core to record the execution time with precision of 0.01 s. Since the main program of an IoT application is usually a big loop, in the experiments we measure the execution time of 1000 loops

 TABLE II

 TOTAL EXECUTION TIME (IN SECOND) OF 1000 LOOPS AND OVERHEADS

Application	# of cleanings	w/o fASLR	with fASLR	Overhead

its semantics. ILR [26] is an instruction-based randomization scheme which relocates every instruction thereby achieving high randomization entropy.

Although fine-grained ASLR is effective in mitigating a single-memory disclosure attack, Snow et al. [3] found that multiple memory disclosures are promising in bypassing finegrained randomization techniques. Motivated by this observation, they introduce an attack framework which bypasses finegrained randomization via just-in-time code reuse (JIT-ROP). With the knowledge of a single-memory disclosure, the framework is able to excavate memory contents of multiple memory pages at runtime, search and assemble gadgets on-the-fly, and then launch CRA. Accordingly a fine-grained randomization approach named Isomeron [27] is proposed as the countermeasure to JIT-ROP attacks. Combining fine-grained ASLR with execution path randomization, Isomeron makes any gadgets unpredictable. Specifically, it generates diversified application code using fine-grained ASLR, and loads both original code and diversified code to the virtual address space at runtime. During execution, a coin-flip decision is made upon each function call to select the destination from either original or diversified code. Related research has been performed to overcome newly emerging CRAs and meet increasing compatibility requirements [28]–[31].

Shi *et al.* [5] leveraged the TrustZone-M hardware extension to enable a function-level ASLR scheme for ARM-based MCUs. The proposed system loads the NS code to NS RAM and periodically reordering all functions at runtime. Compared with our work, this scheme loads the whole application code to RAM. Instead of loading the whole NS app code, our mechanism—fASLR—only loads functions in use and cleans up finished functions from RAM at runtime. fASLR requires smaller RAM and achieves a larger randomization entropy for resource-constrained IoT devices. Shi *et al.* [5] rewrote binaries of the NS code offline and introduces a code size overhead of about 10%–15%, while fASLR has a code size overhead below 5%.

X. CONCLUSION

In this article, we propose fASLR for runtime software security of resource-constrained IoT devices, particularly those based on microcontrollers. fASLR leverages hardware-based security provided by the TrustZone-M technique as the trust anchor. It uses MPU and prevents direct code execution of the application image in the NSW flash. Instead, it traps control flow in an exception handler and relocates functions to be executed to a randomly selected location within the RAM. A memory management strategy was designed for allocating and cleaning up functions in the RR. We also optimized the baseline function cleaning scheme to largely decrease runtime overhead. fASLR is user friendly and only requires a user compiling the app with specific flags. We formally prove that fASLR will not affect the correctness of the NS app execution. We implemented fASLR with a TrustZone-M-enabled MCU-SAM L11. fASLR achieves high randomization entropy with acceptable overheads. We will release fASLR to GitHub for broad adoption and refine the implementation to further reduce the overhead.

REFERENCES

- [1] S. Chen, J. Xu, N. Nakka, Z. Kalbarczyk, and R. K. Iyer, "Defeating memory corruption attacks via pointer taintedness detection," in *Proc. Int. Conf. Depend. Syst. Netw. (DSN)*, Jul. 2005, pp. 378–387. [Online]. Available: https://doi.org/10.1109/DSN.2005.36
- [2] T. K. Bletsch, X. Jiang, and V. W. Freeh, "Mitigating code-reuse attacks with control-flow locking," in *Proc. 27th Annu. Comput. Security Appl. Conf. (ACSAC)*, Orlando, FL, USA, Dec. 2011, pp. 353–362. [Online]. Available: https://doi.org/10.1145/2076732.2076783
- [3] K. Z. Snow, F. Monrose, L. Davi, A. Dmitrienko, C. Liebchen, and A.-R. Sadeghi, "Just-in-time code reuse: On the effectiveness of finegrained address space layout randomization," in *Proc. IEEE Symp. Security Privacy (SP)*, Berkeley, CA, USA, May 2013, pp. 574–588. [Online]. Available: https://doi.org/10.1109/SP.2013.45
- [4] "TrustZone for cortex-M." ARM. [Online]. Available: https:// www.arm.com/why-arm/technologies/trustzone-for-cortex-m (Accessed: Jun. 18, 2022).
- [5] J. Shi, L. Guan, W. Li, D. Zhang, P. Chen, and P. Chen, "HARM: Hardware-assisted continuous re-randomization for microcontrollers," in *Proc. IEEE Eur. Symp. Security Privacy (EuroS P)*, 2022, pp. 520–536.
- [6] X. Shao, L. Luo, Z. Ling, H. Yan, Y. Wei, and X. Fu, "fASLR: Functionbased ASLR for resource-constrained IoT systems," in *Proc. ESORICS*, 2022.
- [7] "freeRTOS—Market Leading RTOS (Real Time Operating System for Microcontrollers)." [Online]. Available: https://www.freertos.org/ (Accessed: Jun. 18, 2022).
- [8] "ARMv8-M Fault Handling and Detection." ARM. [Online]. Available: https://developer.arm.com/documentation/100691/0200/Fault-exceptions (Accessed: Jun. 18, 2022).
- J. Yiu, "Chapter 2—Getting started with cortex-M programming," in Definitive Guide to Arm[®] Cortex[®]-M23 and Cortex-M33 Processors, J. Yiu, Ed. Cambridge, MA, USA: Newnes, 2021, pp. 19–51. [Online]. Available: https://www.sciencedirect.com/science/article/pii/ B9780128207352000020
- [10] S. M. Hejazi, C. Talhi, and M. Debbabi, "Extraction of forensically sensitive information from windows physical memory," *Digit. Investig.*, vol. 6, pp. S121–S131, Sep. 2009. [Online]. Available: https:/ /www.sciencedirect.com/science/article/pii/S1742287609000474
- [11] A. Follner, A. Bartel, and E. Bodden, "Analyzing the gadgets," in Proc. Int. Symp. Eng. Secure Softw. Syst., 2016, pp. 155–172.
- [12] M. D. Brown and S. Pande, "Is less really more? why reducing code reuse gadget counts via software debloating doesn't necessarily indicate improved security," 2019, arXiv:1902.10880.
- [13] "CPU Benchmark—MCU Benchmark—CoreMark." Embedded Microprocessor Benchmark Consortium. [Online]. Available: https:// www.eembc.org/coremark/ (Accessed: Jun. 18, 2022).
- [14] H. Quinn. "Microcontroller Benchmark Codes for Radiation Testing." Los Alamos National Security. [Online]. Available: https://github.com/ lanl/benchmark_codes (Accessed: Jun. 18, 2022).
- [15] J. Pallister, S. Hollis, and J. Bennett, "BEEBS: Open benchmarks for energy measurements on embedded platforms," 2013, arXiv:1308.5174.
- [16] "ATMEL Start." Microchip. [Online]. Available: https://start.atmel.com/ (Accessed: Jun. 18, 2022).
- [17] "STM32L562E-DK—Discovery Kit With STM32L562QE MCU." STMicroelectronics. [Online]. Available: https://www.st.com/en/ evaluation-tools/stm32l562e-dk.html (Accessed: Jun. 18, 2022).
- [18] SWIAT. "On the Effectiveness of DEP and ASLR." Microsoft Security Response Center. 2010. [Online]. Available: https://msrcblog.microsoft.com/2010/12/08/on-the-effectiveness-of-dep-and-aslr/
- [19] "ARM11 MPCore Processor Technical Reference Manual." ARM. [Online]. Available: https://developer.arm.com/documentation/ddi0360/ f/memory-management-unit/memory-access-control/execute-never-bits (Accessed: Jun. 18, 2022).
- [20] "Apply Mitigations to Help Prevent Attacks Through Vulnerabilities." Microsoft Docs. 2021. [Online]. Available: https://docs.microsoft. com/en-us/microsoft-365/security/defender-endpoint/exploit-protection? view=o365-worldwide
- [21] T. Nyman, J.-E. Ekberg, L. Davi, and N. Asokan, "CFI CaRE: Hardwaresupported call and return enforcement for commercial microcontrollers," in *Proc. Int. Symp. Res. Attacks Intrusions Defenses*, 2017, pp. 259–284.
- [22] P. Larsen, A. Homescu, S. Brunthaler, and M. Franz, "SoK: Automated software diversity," in *Proc. IEEE Symp. Security Privacy*, 2014, pp. 276–291.

- [23] C. Kil, J. Jun, C. Bookholt, J. Xu, and P. Ning, "Address space layout permutation (ASLP): Towards fine-grained randomization of commodity software," in *Proc. 22nd Annu. Comput. Security Appl. Conf. (ACSAC)*, 2006, pp. 339–348.
- [24] R. Wartell, V. Mohan, K. W. Hamlen, and Z. Lin, "Binary stirring: Selfrandomizing instruction addresses of legacy x86 binary code," in *Proc. ACM Conf. Comput. Commun. Security*, 2012, pp. 157–168.
- [25] L. V. Davi, A. Dmitrienko, S. Nürnberger, and A.-R. Sadeghi, "Gadge me if you can: Secure and efficient ad-hoc instruction-level randomization for x86 and ARM," in *Proc. 8th ACM SIGSAC Symp. Inf. Comput. Commun. Security*, 2013, pp. 299–310.
- [26] J. Hiser, A. Nguyen-Tuong, M. Co, M. Hall, and J. W. Davidson, "ILR: Where'd my gadgets go?" in *Proc. IEEE Symp. Security Privacy*, 2012, pp. 571–585.
- [27] L. Davi, C. Liebchen, A.-R. Sadeghi, K. Z. Snow, and F. Monrose, "Isomeron: Code randomization resilient to (just-in-time) return-oriented programming," in *Proc. NDSS*, 2015, pp. 1–15.
- [28] H. Koo, Y. Chen, L. Lu, V. P. Kemerlis, and M. Polychronakis, "Compiler-assisted code randomization," in *Proc. IEEE Symp. Security Privacy (SP)*, 2018, pp. 461–477.
- [29] F. Xuewei, W. Dongxia, L. Zhechao, K. Xiaohui, and Z. Gang, "Enhancing randomization entropy of x86-64 code while preserving semantic consistency," in *Proc. IEEE 19th Int. Conf. Trust Security Privacy Comput. Commun. (TrustCom)*, 2020, pp. 1–12.
- [30] S. Priyadarshan, H. Nguyen, and R. Sekar, "Practical fine-grained binary code randomization," in *Proc. Annu. Comput. Security Appl. Conf.*, 2020, pp. 401–414.
- [31] X. Wang, S. Yeoh, R. Lyerly, P. Olivier, S.-H. Kim, and B. Ravindran, "A framework for software diversification with ISA heterogeneity," in *Proc. 23rd Int. Symp. Res. Attacks Intrusions Defenses (RAID)*, 2020, pp. 427–442.



Zhen Ling (Member, IEEE) received the B.S. degree from Nanjing Institute of Technology, Nanjing, China, in 2005, and the Ph.D. degree in computer science from Southeast University, Nanjing, in 2014.

He is a Professor with the School of Computer Science and Engineering, Southeast University. His research interests include network security, privacy, and Internet of Things.

Prof. Ling won the ACM China Doctoral Dissertation Award in 2014 and the China Computer Federation Doctoral Dissertation Award in 2015.

Huaiyu Yan received the B.S. degree in software engineering from Southeast University, Nanjing, China, in 2019, where he is currently pursuing the Ph.D. degree in computer science and engineering.

His current research interests include Internet of Things and privacy and security.



Yumeng Wei is currently pursuing the B.S. degree in cyberspace security with Southeast University, Nanjing, China.

Her research interests include software and network security of Internet of Things devices.



Lan Luo received the B.S. degree in electrical engineering from the Civil Aviation University of China, Tianjin, China, in 2015, and the M.S. degree in computer engineering and the Ph.D. degree in computer science with the University of Central Florida, Orlando, FL, USA, in 2018 and 2022, respectively. Her research interests mainly cover security and

Her research interests mainly cover security and privacy of Internet of Things, security of embedded system, network and software security, and trustworthy computing.



Xinhui Shao received the B.S. degree in communication engineering from Shanghai University, Shanghai, China, in 2019. He is currently pursuing the master's degree in cyber science and engineering with Southeast University, Nanjing, China.

His current research interests include Internet of Things and privacy and security.



Xinwen Fu (Senior Member, IEEE) received the B.S. degree in electrical engineering from Xi'an Jiaotong University, Xi'an, China, in 1995, the M.S. degree in electrical engineering from the University of Science and Technology of China, Hefei, China, in 1998, and the Ph.D. degree in computer engineering from Texas A&M University, College Station, TX, USA, in 2005.

He is a Professor with the Department of Computer Science, University of Massachusetts Lowell, Lowell, MA, USA. He was a tenured

Associate Professor with the Department of Computer Science, University of Central Florida, Orlando, FL, USA. He has published at prestigious conferences, including the four top computer security conferences (Oakland, CCS, USENIX Security, and NDSS), and journals, such as ACM/IEEE TRANSACTIONS ON NETWORKING and IEEE TRANSACTIONS ON DEPENDABLE AND SECURE COMPUTING. His current research interests are in computer and network security and privacy.

Dr. Fu spoke at various technical security conferences including Black Hat.