

SIC²: Securing Microcontroller Based IoT Devices with Low-cost Crypto Coprocessors

Bryan Pearson*, Cliff Zou*, Yue Zhang^{†§}, Zhen Ling[‡], Xinwen Fu[†]

*Dept. of Computer Science, University of Central Florida, USA

Email: bpearson@knights.ucf.edu, czou@cs.ucf.edu

[‡]School of Computer Science and Engineering, Southeast University

Email: zhenling@seu.edu.cn

[†]Dept. of Computer Science, University of Massachusetts Lowell, Lowell, MA, USA

Email: xinwen_fu@uml.edu

[§]Dept. of Computer Science, Jinan University

Email: zyueinfosec@gmail.com

Abstract—In this paper, we explore the use of microcontrollers (MCUs) and crypto coprocessors to secure IoT applications, and we show how developers may implement a low-cost platform that provides security to users and protects private keys against software attacks. We first demonstrate the plausibility of format string attacks on the ESP32, a popular MCU from Espressif that uses the Harvard architecture. The format string attacks can be used to remotely steal private keys hard-coded in the firmware. We then present a framework termed *SIC²* (Securing IoT with Crypto Coprocessors), for secure key provisioning that protects end users’ private keys from both software attacks and untrustworthy manufacturers. As a proof of concept, we pair the ESP32 with the low-cost ATECC608A cryptographic coprocessor by Microchip and connect to Amazon Web Services (AWS) and Amazon Elastic Container Service (EC2) using a hardware-protected private key, which provides the security features of TLS communication including authentication, encryption and integrity. We have developed a prototype and performed extensive experiments that show that the ATECC608A crypto chip may significantly reduce the TLS handshake time by as much as 82% with the remote server, and it may lower the total energy consumption of the system by up to 70%. Our results indicate that securing IoT with crypto coprocessors is a practicable solution for low-cost MCU based IoT devices.

I. INTRODUCTION

The popularity of Internet of Things (IoT) has raised grave security and privacy concerns. There is a broad attack surface against IoT, including vulnerabilities and issues in hardware, firmware/operating system, application software, networking and data. For example, hackers can force autonomous vehicles to crash [21] and may also steal credentials from consumer and medical products [2]. Botnets such as Mirai [1] and Reaper [7] exposed vulnerable networks and compromised millions of devices.

IoT device manufacturers have been advancing the hardware to secure IoT devices. One of the pioneers is Espressif Systems, which produces the popular ESP8266 and ESP32 chips and claimed a shipment of 100 million of both chips in January 2020 [26]. Particularly, ESP32 has abundant hardware security features including secure boot and flash encryption [28]. However, we have found potential threats

against ESP32 by using two practical software attacks, named Same Subroutine Attack and Cross Subroutine Attack. In Same Subroutine Attack, the vulnerable instruction (i.e., `printf(.)` or `sprintf(.)`) and the victim’s secret data are co-located in the same subroutine, and an attacker may steal this secret data via the attack. However, Cross Subroutine Attack is more powerful, where an attacker may extract sensitive information even if the vulnerable function and the victim code fragment are located in different subroutines. We demonstrate our attacks with a proof of concept web server, showing that an attacker may deploy the attacks remotely through the Internet. Our attacks significantly undermine the security of ESP32, and the principles may apply to other IoT chips.

To defeat software attacks, we explore the use of low-cost cryptographic coprocessors (costing less than \$1) to secure low-cost IoT devices based on microcontrollers (MCUs). With a cryptographic coprocessor chip that can serve as the root of trust, private keys may never leave the chip, and cryptographic operations over data from the main MCU are performed inside the chip. We present a secure key provisioning solution, denoted as *SIC²* or Securing IoT devices with Cryptographic Coprocessors, that stores private keys inside of a cryptographic coprocessor. Our *SIC²* protects keys from malicious personnel within the semiconductor manufacturing line as well as cyber attacks against those manufacturers [14]. We implement a proof of concept by pairing the ESP32 with the ATECC608A [9] crypto coprocessor (\$0.53 at Microchip), which can provide mutual authentication, encryption and integrity to a network.

Our major contributions can be summarized as follows:

- 1) We show that popular MCUs such as ESP32 can be compromised by multiple software attacks. Private keys can be leaked remotely.
- 2) We propose *SIC²*, a systematic solution for manufacturers to securely write private keys into cryptographic coprocessors to secure IoT devices. We use ESP32 as an example, pairing the MCU with a new cryptographic coprocessor ECC608A. We offer design and

implementation criteria for developers.

- 3) We perform extensive experiments to validate the speed performance and energy consumption of SIC^2 . Our results show that connecting to a cloud server such as Amazon EC2 can reduce the overall TLS handshake time by 82% and energy consumption by up to 70%.

The rest of this paper is organized as follows. In Section II, we provide the background of the ESP32 MCU and its processor. In Section III, we present novel format string attacks against the ESP32 which compromise private keys stored on the device. In Section IV, we introduce SIC^2 and how manufacturers may securely write private keys into cryptographic coprocessors. A proof of concept of SIC^2 is discussed in Section V which combines the ESP32 with the ECC608A. We evaluate the ECC608 overhead and network performance in Section VI. Section VII discusses some related works, and Section VIII concludes the paper.

II. BACKGROUND

In this section, we discuss the system design of the ESP32 and the architecture of the Xtensa processor, which is used by the ESP32.

A. ESP32 System Design

The ESP32 is a popular IoT MCU used by millions of consumers around the world [26]. The ESP32 contains a network stack that supports WiFi, Bluetooth, and Bluetooth Low Energy (BLE) capabilities for a variety of IoT applications. It exposes Universal Asynchronous Receiver/Transmitter (UART) and Joint Test Action Group (JTAG) external debugging ports. UART communication allows users to monitor console output, upload new firmware to the chip, dump flash contents, and modify security settings on the chip. JTAG allows for complete debugging of the ESP32, including reading and modifying the entire firmware, bootloader contents, CPU registers and SRAM contents on the ESP32. Espressif has ported GDB to recognize the Xtensa architecture.

The ESP32 contains a 1kB block of secure eFuse memory. These memory contents are accessible only to the hardware, and once an eFuse value is set, it is irreversible. The eFuse memory controls access to the communication and debugging ports. Another feature of the eFuse is the secure storage of a 256-bit flash encryption key and a 256-bit secure boot key. With flash encryption enabled, the ESP32 can use the flash encryption key with AES cipher block chaining (CBC) mode to decrypt data and instructions before being processed by the CPU. With secure boot, the ROM will calculate an AES-based SHA digest to validate the integrity of the bootloader, which in turn may validate the firmware.

B. Xtensa Processor Architecture

In this section, we discuss architecture details about Xtensa LX6, a 32-bit microprocessor from Tensilica [12]. ESP32 contains 2 Xtensa processors. We first provide some basic information about the Xtensa processor. We then discuss details about the register file and how the ESP32 can access register contents at runtime.

1) *Architecture Details:* Xtensa implements a modified Harvard architecture [25], with the main memory separated into instruction SRAM and data SRAM. The processor is programmable to allow manufacturers to modify instructions, the register file size, cache size, memory width, and make various other enhancements. Tensilica provides tools for mapping any configuration to the physical hardware.



Fig. 1: Overview of the ESP32 register file. A subroutine only has access to the registers contained within the register window.

2) *The Register Window:* The ESP32's register file contains 64 general-purpose registers. The Xtensa architecture implements a feature called **register window** which allows a subroutine to only access to 16 general-purpose registers at a time. In the register file, registers are labeled AR0, AR1, AR2, etc. The register window allocates a contiguous block of registers within the register file; for example, a subroutine may only have access to AR16, AR17, AR18, and so forth, up to AR31. When a subroutine *Sub1* calls some other subroutine *Sub2*, the register window "increments" its position in the register file, meaning new registers become available while old registers become inaccessible. The register window can increment by either 4, 8, or 12 registers. When *Sub2* returns, the register window reverts or "decrements" to the original position, allowing *Sub1* to access the same registers as before.

Figure 1 provides further details. Consider *Sub1*, whose register window is defined for the range AR16 to AR31. Now when *Sub1* calls *Sub2*, the register window increments by 8 registers such that *Sub2* can now access registers in the range AR24 to AR39, while registers AR16 to AR23 are no longer accessible. Similarly, when *Sub2* calls *Sub3*, the register window increments by 8 registers again such that *Sub3* can access registers in the range AR32 to AR47. On each return—from *Sub3* to *Sub2* and from *Sub2* to *Sub1*—the register window will decrement by 8 registers and allow each respective subroutine to recover the contents of its registers.

In the case where a register window attempts to allocate registers that already belong to a parent subroutine *SubP*, the CPU will initiate a *window overflow exception*. In this scenario, the CPU will dump the contents of registers into memory and allow the new subroutine to access those regis-

ters. When the program returns back to *SubP*, it will restore the register contents from memory back into the registers. In this way, register contents are never lost, even when the registers themselves must be shared among subroutines.

III. NOVEL ATTACKS AGAINST ESP32

In this section, we present several new attack proof-of-concepts against the ESP32. We have successfully launched two format string attacks on the ESP32, named the **Same Subroutine Attack** and the **Cross Subroutine Attack**. We begin by explaining the threat model of these attacks. Then we discuss implementation details, before showing a proof of concept for a remote format string attack. Additional proof-of-concepts, including a "Serial-to-TCP" attack, are discussed in Appendix D.

A. Threat Model

In the following attacks, we assume that the ESP32 may expose some kind of communication channel to the user, such as HTTP or MQTT. We also assume that the ESP32 stores some secret data in its firmware. The adversary may be local or remote. A local adversary can physically access the device and view serial output directly. In the local attack, we assume that the adversary can access the UART interface on the ESP32. A remote adversary can read the data transmitted by the ESP32 over the communication channel. We show that the proposed attacks can work both locally and remotely. Finally, we assume that the firmware contains some programming flaw, which is reasonable due to the abundance of software vulnerability types in C [31] [32] [30] [24].

B. Attack Overview

Format string vulnerabilities [30] arise when formatting functions fail to validate a user's input format. An example of such a function is *printf()*, which accepts format string characters as its input. Typically, if the program were to execute an instruction such as *printf("%s", name)*, it would simply print the contents of *name*. However, if *name* is not provided to the function, the program will print the contents of a different memory location, which may leak sensitive data. Every format string character passed to the format function will fetch the value of the next consecutive memory address and cast it accordingly. On the ESP32, which has a 4-byte address width, this means every format string character fetches the next 4 bytes in memory.

Based on our experiments on the ESP32, we found that when no input parameters are provided to *printf()*, it will begin by fetching the last five registers in the subroutine's register window. Afterwards, it will fetch the value at the stack pointer (SP), then the value at SP + 4, then SP + 8, and so forth. This means that on the ESP32, at least 6 format string characters are required to begin accessing memory contents. In this way, the format string attack may be used on the ESP32 to leak arbitrary data from memory.

C. Format String Attacks

1) **Same Subroutine Attack**: In the **Same Subroutine Attack**, the format string instruction and the private data exist within the same subroutine. We begin by discussing the setup of this subroutine. We then describe the details of the registers and memory. Finally, we show how an adversary may exploit this program and obtain the private data.

Listing 1: Same Subroutine Attack psuedocode.

```
void app_main() {
    char tmp[16] = "PRIVATE KEY";
    char* params = malloc(128);
    accept_user_input(&params);
    printf(params); }
```

As shown by Listing 1, the program defines a local variable called *tmp* in a subroutine called *app_main*. In our example, *tmp* is a 16-byte char array set to the string "PRIVATE KEY". *printf()* will print some arbitrary input that is provided by the user in *accept_user_input()*.

We used JTAG debugging on the ESP32 to determine details about this program. Communication with JTAG requires the addition of OpenOCD, an open-source software project that can communicate with the JTAG interface [23]. We attached a GDB client to the OpenOCD session in order to debug our application.

From JTAG debugging, we determined the following details. The stack pointer address of *app_main* is 0x3ffb4ee0. On the ESP32, local char arrays are always defined starting at the stack pointer address, so *tmp* is defined from 0x3ffb4ee0 to 0x3ffb4eef. The register window in the subroutine is defined between AR16 and AR31. The contents of the last five registers in the register window (namely, AR27 to AR31) are 0x8001f880, 0x6ff1ff8, 0x0, 0x3ffaffe0, and 0x3ffb6840.

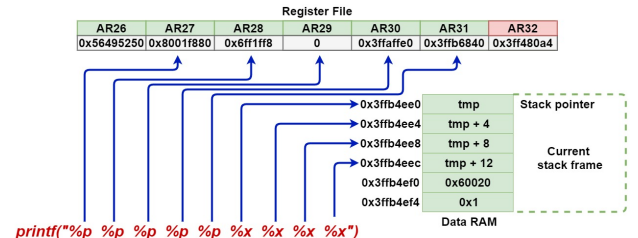


Fig. 2: Overview of the Same Subroutine Attack. The arrows show which addresses *printf()* will access.

To perform the attack, the user must provide the following format string as input to *printf()*: "%p %p %p %p %p %x %x %x %x". The first five characters print the contents of AR27, AR28, AR29, AR30, and AR31, while the last four characters print the contents of *tmp*. The attack behavior is illustrated in Figure 2. The output to UART is shown below:

```
0x8001f880 0x6ff1ff8 0x0 0x3ffaffe0
0x3ffb6840 56495250 20455441 59454b
0
```

As shown above, the first five values correspond to the register contents of AR27 through AR31. The next 16 bytes correspond to the stack contents, beginning with the stack pointer. Recall that *tmp* is written at the stack pointer address.

Since the bus architecture of the ESP32 is little-endian, the bytes must be reversed to recover the original data. For example, the value 56495250 must be changed to 50524956. After doing this for all values, the user can obtain the desired value 50524956415445204b45590. A hex-to-ascii converter shall reveal the contents of this data to be "PRIVATE KEY".

2) **Cross Subroutine Attack:** In the **Cross Subroutine Attack**, the format string instruction and the private data are located in different subroutines. This attack is much more powerful than the Same Subroutine Attack, since it can steal data from any previous subroutine in the call stack. Again, we begin by discussing the setup requirements of the program, followed by the details of the program including memory and register contents. We conclude by showing the exploit and how the private data may be recovered.

Listing 2 shows that the format string function and the private data are located in different subroutines. we have a local variable *tmp* defined in *app_main*, but we also have two new subroutines, *sub1* and *sub2*. In the expected program flow, *app_main* calls *sub1*, which calls *sub2*, which calls the vulnerable *printf* function. The attack will leverage the behavior of the window overflow exception in the ESP32, where register contents are dumped to memory when the program transfers control to a new subroutine. Namely, the address of *app_main*'s stack pointer will dump to memory when the program executes *sub2*, and the attacker can use the character "%s" at this location to recover the stack pointer address, cast it as a string, and print its value.

Listing 2: Cross Subroutine Attack psuedocode.

```
void sub2(char* x) { printf(x); }
void sub1(){
    char* params = malloc(128);
    accept_user_input(&params);
    sub2(params); }
void app_main() {
    char tmp[16] = "PRIVATE KEY";
    sub1(); }
```

Again, we used JTAG to debug the program and discovered the following information. First, the stack pointer of *app_main*, *sub1*, and *sub2* are 0x3ffb4ee0, 0x3ffb4ec0, and 0x3ffb4ea0, respectively. As *tmp* is a local buffer defined in *app_main*, *tmp*'s address is also 0x3ffb4ee0. The ESP32's application startup sequence makes several subroutine calls prior to reaching *app_main*, and our experiments show that the register file has already been exhausted by the time the program reaches *app_main*. The call to *sub1* and *sub2* both shift the register window by 8 registers. Therefore, due to the window overflow exception, 8 registers must be dumped into memory on both calls. The register window for *app_main* is defined from AR16 to AR31. For *sub1*, it is defined from AR24 to AR39. And for *sub2*, it is defined from AR32 to AR47. When jumping to *sub1*, registers AR16 through AR23 are saved to memory; when jumping to *sub2*, registers AR24 through AR31 are saved to memory. The stack pointer address is always stored in the second register of the register window; in the case of *app_main*, AR17 contains the stack pointer value 0x3ffb4ee0. Our experiments revealed that the

second register is always dumped to the memory location that is 12 bytes behind the new stack pointer. In particular, this means that when the program reaches *sub1*, the stack pointer address of *app_main* is saved to 0x3ffb4eb4, exactly 12 bytes behind *sub1*'s stack pointer.

If *printf* can be manipulated to point to 0x3ffb4eb4, the character "%s" will cast this address as a char pointer and print its value accordingly. This will cause the contents of *tmp* to be leaked. However, as noted above, the stack pointer address of *sub1* is 0x3ffb4ec0, while the stack pointer address of *sub2* is 0x3ffb4ea0. This means that the format string attack cannot be used in *sub1*, because the format string pointer can only be moved forward in memory starting from the stack pointer; it cannot be moved backward. Fortunately, *sub2*'s stack pointer precedes 0x3ffb4eb4, so the format string attack is feasible in this subroutine.

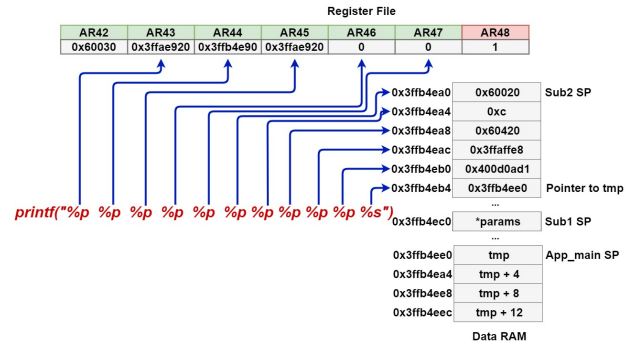


Fig. 3: Overview of the Cross Subroutine Attack.

To perform the attack, the user must provide the following format string as input to *printf*(): "%p %p %p %p %p %p %p %p %p %p %p %s". Note that *sub2*'s stack pointer precedes 0x3ffb4eb4 by 20 bytes. The first 5 characters print the registers AR43 through AR47, while the next 5 print the first 20 bytes of the stack frame, specifically the contents in the range 0x3ffb4ea0 to 0x3ffb4eb3. The final character, "%s", will cast 0x3ffb4eb4 as a char pointer and print its value. The attack behavior is illustrated in Figure 3. The output of this attack is shown below:

```
0x3ffae920 0x3ffb4e90 0x3ffae920
0x0 0x0 0x60020 0xc 0x60420
0x3ffa8 0x400d0ad1 PRIVATE_KEY
```

Unlike the Same Subroutine Attack, this payload prints the memory contents verbatim as long as it can be cast as a string. This can lead to faster detection of sensitive data. One caveat is that this can lead to unexpected errors if a string does not contain the null terminator, which will cause the program to read memory contents indefinitely and potentially crash.

D. Proof of Concept Remote Format String Attack

We now demonstrate how format string attacks can be performed remotely to steal an IoT device's secrets. An application that exposes a communication channel such as a web server interface may fall victim to the format string attack, which can leak secrets to remote attackers. We have written two vulnerable programs using the Arduino

IDE, an alternative to the ESP-IDF development platform provided by Espressif. To serve a web server, ESP32 uses the WebServer library written for the Arduino platform. This library allows a web server to process HTTP requests from the client and send responses back. The ESP32 application contains a format string vulnerability based on the *sprintf()* function in C. The *sprintf()* function sends the formatted output to a string rather than stdout. The expected syntax for this instruction is *sprintf(buf, "%s", param)*, where *buf* is a string and *param* is formatted as a string before being sent to *buf*. However, the instruction *sprintf(buf, param)* is vulnerable to the format string attack, because *param* is now treated as the format parameter and can lead to memory leakage if controlled by an attacker. If *buf* is passed remotely to an adversary, he can observe the output of the format string attack. Using this knowledge, we show how the Same Subroutine Attack and Cross Subroutine Attack can occur remotely.

1) *Remote Same Subroutine Attack*: This attack exploits the Same Subroutine Attack through the *sprintf()* vulnerability to leak a private key on a web server. Listing 3 provides pseudocode, which starts a web server at port 80. The *handleReq()* subroutine is called whenever a user visits the root index of the server. This subroutine stores a key and parses any HTTP GET request sent from the client, calls *sprintf()* to format the request, and passes the formatted request back to the client in an HTTP response.

Listing 3: Remote Same Subroutine Attack pseudocode.

```
WebServer server(80);
void handleReq() {
    char key[32] = "THIS IS A PRIVATE KEY";
    char* res = malloc(128);
    char* param = server.arg(0).c_str();
    sprintf(res, param);
    server.send(200, "text/plain", res);
}
void setup() { server.on("/", handleReq); }
void loop() { server.handleClient(); }
```

To send a payload, the adversary can use a web browser to send the following GET request to the ESP32:

```
http://[IP addr]/?h=%25x+%25x+%25x+
%25x+%25x+%25x+%25x+%25x+%25x+%25x+
%25x+%25x+%25x
```

The server will receive the format string and parse it during the *sprintf()* instruction, which will leak the contents of the private key into the *res* buffer. The server will send the buffer back to the adversary to read in the HTTP response. The attacker can then derive the key from the payload by decoding to ascii.

2) *Remote Cross Subroutine Attack*: This attack exploits the Cross Subroutine Attack through the *sprintf()* vulnerability. Similar to the *printf()*-based Cross Subroutine Attack, the adversary will pass a format string containing the "%s" character to cast the private key as a string. The web server will pass this output to a buffer that is sent back to the adversary in an HTTP response. The full details of this attack are available in Appendix D.

IV. SIC²: SECURING IOT WITH CRYPTO COPROCESSORS

In this section, we discuss the need of crypto coprocessors for IoT devices and present a secure key provisioning framework. Then we provide a security analysis of the framework.

A. Need of Crypto Co-processors

From our discussion in Section III, MCUs with secure boot can be compromised and leak cryptographic keys if these keys have no hardware protection. The TrustZone technology has been integrated into Arm Cortex-M processors, denoted as TrustZone-M. However, TrustZone-M can be compromised too [19]. If an application in a MCU directly accesses cryptographic keys for cryptographic functionalities, once the MCU system is compromised, the cryptographic keys will leak. Therefore, a crypto coprocessor chip is an ideal solution. The application feeds data to the crypto coprocessor, which stores the keys, performs cryptographic functionalities inside the chip and returns the results to the application in the MCU.

We have examined over 40 MCUs and a number of IoT development boards and solutions. Only Microsoft's Azure Sphere [20] and TI's CC3220 and CC3100MOD have integrated crypto coprocessors with the MCUs. Fortunately, there are two standalone crypto coprocessor modules, Microchip's ATECC608/ATECC508 (around \$0.53/unit) and NXP's SE050 (around \$0.97/unit). Only a few development boards have begun to use these crypto coprocessor modules, including Microchip's SAM L11 Xplained Pro Evaluation Kit and Arduino NANO 33 IOT. Our full dataset is provided in Appendix A.

B. Secure Key Provisioning

We introduce our secure key provisioning model, which allows an IoT manufacturer to adopt low-cost crypto coprocessors without leaking secret keys written into the crypto coprocessors. Manufacturers will defer the provisioning of private keys and certificates to a **secure facility**, which is separated from the rest of the manufacturing process and responsible for storing data inside the crypto chips. Even this secure facility cannot access private keys, which are internally generated by the crypto coprocessor.

Secure key provisioning is a grand challenge while incorporating a crypto coprocessor into an IoT system. Without secure provisioning, private keys may be leaked by malicious personnel within the manufacturing line or by supply-chain attacks [14]. An ideal IoT solution is that each IoT device has at least one unique private key (in terms of public key cryptography) along with a certificate stored in the secure storage of the crypto coprocessor, and the public key associated with the crypto coprocessor can be safely derived by the party who wants it. To solve this key provisioning problem, we have to answer questions such as: who will inject a private key into the crypto coprocessor? And when? We provide a novel framework considering the entire development cycle of the IoT system.

Our secure key provisioning framework is shown in Figure 4. It is composed of five main entities. The **factory** is a

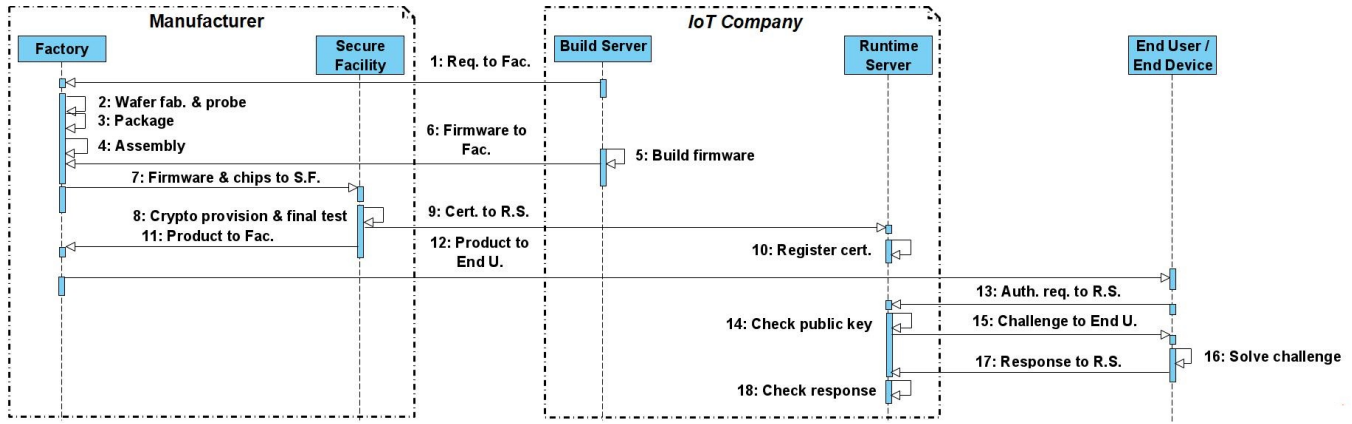


Fig. 4: Secure key provisioning framework.

generic concept that will represent the complete semiconductor manufacturing line, which can be widely varied. This includes the fabrication, packaging, assembly, and testing of the hardware. The factory will manufacture crypto chips and IoT devices. Additionally, end users can purchase their IoT products from the factory. The **secure facility** will receive crypto chips from the factory and provision them with private keys and certificates. The secure facility will also distribute these certificates to the runtime server. The **build server** creates the firmware for the end device. The **runtime server** serves as the application server and authenticates the end device's public key and certificate. Finally, the **end user / end device** is the final IoT product / the owner of the final IoT product. The factory and secure facility are part of the generic **Manufacturer** group, while the build server and runtime server are specific to the **IoT Company**.

1) *Manufacturing Phase (Steps 1-4)*: In this phase, the manufacturing line produces the hardware of the chips according to the specification by the IoT company. The build server will send a manufacturing request to the factory, including hardware requirements and the identity of the runtime server. The factory follows this request to manufacture and assemble the IoT device. This includes four key steps. **Wafer fabrication** constructs the silicon die to connect the electrical components together. **Wafer probing** performs electrical tests on the silicon chip. **Packaging** packages the die (i.e., block of semiconducting material) to protect the electrical components from damage. And **assembly** refers to the production of printed circuit boards (PCBs) and assembling the modules and chips onto the PCB. Assembly may occur either before or after the key provisioning phase.

2) *Key Provisioning Phase (Steps 5-10)*: In this phase, the secure facility performs the key provisioning process on the crypto coprocessor and generates the device certificates. After developing the product firmware, the build server sends it to the factory, who forwards it to the secure facility. The factory also notifies the secure facility about the runtime server's identity. Then the secure facility provisions each crypto chip to internally generate a private key. Additionally, the secure facility will generate and store a unique device certificate into the device. The certificate identifier, e.g., its Common Name, should be unique to each certificate; for

instance, it can be derived from the identity of the crypto chip, such as a serial number. Next, the secure facility will configure the chip such that its public key and certificate are readable and its private key is locked from read/write access. Finally, the secure facility will upload the firmware to the chip and perform some final testing to ensure that the crypto coprocessor has been correctly provisioned. The secure facility distributes the device certificate to the runtime server, who shall save these certificates to a registry.

3) *Device Authentication Phase (Steps 11-18)*: In this phase, the end user obtains the finished product and authenticates it to the runtime server using the key stored on the crypto chip. The end user will order the product from the factory. Then the user turns on the device and sends an authentication request to the runtime server, which includes the public key of the crypto coprocessor. The runtime server searches its certificate registry to ensure the validity of the public key. Then it will initiate a challenge-response procedure to ensure that the end device owns the public key. The end device will use its private key to sign a challenge and prove ownership of the key. Once authentication is complete, the runtime server and end user can proceed with the normal application.

C. Security Analysis

With our defense enabled, all software attacks in this paper will fail because private keys will no longer be stored in the firmware. Based on the framework, it can be seen that the crypto chip is provisioned in a secure environment, and that a malicious user or factory worker can never steal the private key. One issue is that a factory will manufacture many different products, and the runtime server must only accept certificates which belong to its own products. To address this, the runtime server will receive certificates from the secure facility and can know ahead of time which certificates to trust. In this way, the runtime server will reject certificates from devices that were not provisioned by the secure facility. Additionally, the framework can be extended to provision private keys for other devices besides crypto coprocessors. We present a proof of concept with the ESP32 in Section V-C.

V. PROOF OF CONCEPT OF SIC^2

As a proof of concept, we have implemented SIC^2 via the ESP32 and ECC608 to achieve software security. The ECC608 chip will store a 256-bit ECC private key that can serve as the root of trust for many applications, including network security via X.509 certificates and the TLS cryptographic protocol. In the case of a software exploit, the developer does not need to worry that the private key has been compromised, since the key will be stored in the secure ECC608 chip instead of the compromised ESP32 chip. In addition, the ECC608 provides hardware acceleration of cryptographic functions such as ECDH and ECDSA, allowing the ESP32 to authenticate to a network faster. Furthermore, we have combined the ESP32 and ECC608 with the DHT22 temperature and humidity sensor from Adafruit [13]. A prototype of our defense can be found in Figure 5. This project was written in ESP-IDF version 4.0 and is publicly available on Github ¹.

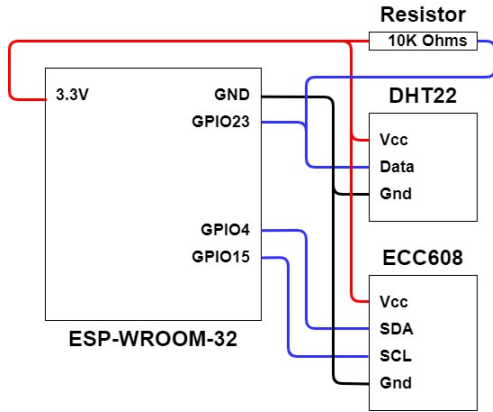


Fig. 5: Schematic of ESP32 with ECC608 and DHT22.

A. ATECC608A Overview

The ECC608 comes packaged in the Small Outline IC (SOIC) format. In the manufacturing line, the SOIC may be directly soldered onto a PCB for maximum area efficiency. Alternatively, a user may solder the SOIC to a socket adapter which can be used on a breadboard. Figure 6 illustrates the pairing of an ESP32-based development board with the ECC608 on a socket adapter.

The ECC608 contains an EEPROM which is capable of storing up to 16 keys, certificates, or user data. Storage regions are organized into *slots*. The slot and its corresponding key may be configured in various ways. Our configuration allows the ECC608 to generate and verify signatures and extract the public key. The private key cannot be read or modified. The ECC608 may also generate a certificate signing request (CSR) from the private key. This is necessary for attaining a valid X.509 certificate. To prevent malicious configuration or overwriting of data, the user should lock the *configuration* and *data* memory zones.

A device can communicate with the ECC608 via the CryptoAuthLib software library [11]. CryptoAuthLib allows an MCU to communicate with the ECC608 via the I^2C protocol

to lock the memory zones and send other commands. The host MCU and ECC608 may also share a mutual input/output secret, which obscures the I^2C traffic by encrypting data with the secret value. This results in a safer I^2C channel.

To achieve network communication, we use *MbedTLS* [8], a lightweight crypto library that implements TLS functions on embedded systems. We have modified this library to out-source private key operations to the ECC608. The most critical of these operations is the signature generation function, which is used to sign a challenge packet from the server and prove ownership of a certificate. We have also added support for signature verification and ECDH establishment, in case the server provides an ECC-based certificate. Altogether, the necessary modifications to MbedTLS are quite minimal, as the majority of the code base remains untouched.

Apart from secure key storage, the ECC608 can serve a WiFi-enabled application in other ways. For instance, the ECC608 provides a secure boot feature that can validate a firmware; this can provide additional security to chips such as Arduino or ESP8266. If the ECC608 stores the device certificate or CA certificate, then TLS performance could potentially increase even further. Finally, each ECC608 contains a 72-bit unique serial number that can be used to identify the chip.

B. Integration with ESP32

To combine the ESP32 with the ECC608, we provide details for a complete hardware and software implementation. The CryptoAuthLib and MbedTLS libraries must be ported correctly to compile within ESP-IDF's build system. We provide implementation details with the DHT22 in the technical report [22].

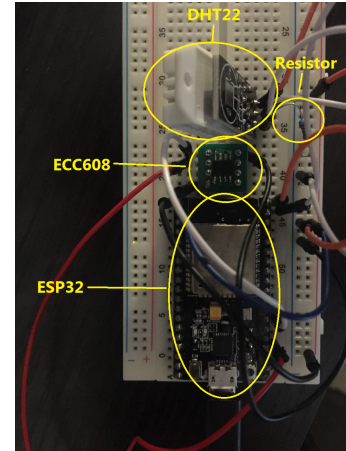


Fig. 6: ESP32 with the ECC608 and DHT22 on a breadboard.

We have paired the crypto chip with a development board that incorporates ESP-WROOM-32 module and 4 MB external flash. To utilize the I^2C interface, we use GPIO ports 15 (SCL) and 4 (SDA) on the ESP32, although other ports such as 21 and 22 can be used. The power supply of the ECC608 connects to the ESP32's 3.3V output pin. We have soldered the ECC608 to a SOIC socket adapter. Figure 6 illustrates our hardware setup on a breadboard.

We have used Atmel Crypto Evaluation Studio (ACES) to set the configuration parameters. ACES is a programming

¹Available at <https://github.com/PBearson/ESP32-With-ECC608>.

software that can communicate with the ECC608 via an external programmer, such as the ATSAM21 board [10].

We have developed a provisioning app that generates an ECC private key in slot 0 and corresponding X.509 CSR. It will also lock the data zone once the private key is set. To port CryptoAuthLib to ESP-IDF, we have cloned the source code from Github and added a "CMakeLists.txt" (a file executed by CMake which describes the build instructions for a project) to the root directory. The file includes the source and header files of this library. The library contains a hardware abstraction layer that specifies communication settings with many devices including the ESP32 over I^2C ; this setting is included as a compile option in "CMakeLists.txt".

In addition, we have developed an app that connects with a remote server via Message Queuing Telemetry Transport (MQTT) over TLS. In our prototype, we connect to an EC2 node where the server and CA certificates have been generated using ECC private keys. In this way, the ECC608 can be used to verify these certificates and generate the session key via ECDH.

Our app integrates the CryptoAuthLib and Espressif's MbedTLS libraries. Like CryptoAuthLib, we write a "CMakeLists.txt" file for MbedTLS that includes the required source files as well as dependencies to CryptoAuthLib. We have modified the ECDSA and ECDH source files included in MbedTLS. We have written alternative functions in these source files which can be enabled or disabled in the port directory, via a configuration file. In ECDSA, we write function overloads for signature generation and signature validation which offload these operations to the ECC608. *atcab_sign* and *atcab_verify_extern* will provide the required operations. In ECDH, we overload the public key generation and shared key generation functions. *atcab_genkey* will generate a key in the temporary key slot, while *atcab_ecdh_tempkey* will establish the shared key.

1) *DHT22 Integration*: The DHT22 is a digital temperature/humidity sensor from Adafruit. The DHT22 can pass readings to an MCU, such as the ESP32, for further processing. As shown in Figure 5, the DHT22 only requires a single data pin as well as voltage and ground. The user should also connect a $10k\Omega$ pull-up resistor between voltage and data to stabilize the input voltage, although some DHT22 breakout boards include an internal pull-up resistor. To compile the sensor for ESP-IDF, the user should include the appropriate library and build it into their project [6]. In the main program file, the user can use the `setDHTgpio(.)` function to set the data pin. Afterwards, the user can use the `getDHT()` method to sample the latest temperature and humidity readings from the DHT22, and the user can use the corresponding `getHumidity()` and `getTemperature()` methods to extract the humidity and temperature values from the reading, respectively. In Listing 4, we provide code that will sample the DHT22 every 3 seconds and print the results to the console:

```
void app_main()
{
    setDHTgpio(23);
    while(1)
    {
        assert(DHT_OK == readDHT());
        float hum = getHumidity();
        float tmp = getTemperature();
        printf("\nHumidity:\t%.1f\n", hum);
        printf("Temperature:\t%.1f\n", tmp);
        vTaskDelay(3000 / portTICK_RATE_MS);
    }
}
```

An output of 3 sample readings is shown below:

```
Humidity:      49.9
Temperature:   27.5

Humidity:      49.9
Temperature:   27.4

Humidity:      49.9
Temperature:   27.4
```

C. Secure Provisioning of the ESP32

The ESP32 can provide flash encryption and secure boot to prevent readout and modification of the firmware. These features rely on two private keys stored in the secure eFuse memory. However, enabling these features presents a challenge due to the security risks involved in key provisioning, as discussed in Section IV.

The details for enabling the ESP32's security features are as follows. 1) **Flash Encryption (FE)**: The programmer should use Espressif's build framework to compile the bootloader to support FE. During the boot sequence, the bootloader will detect FE is supported, and the hardware will generate a key to store in the eFuse. Then, the chip will encrypt the complete flash contents. 2) **Secure Boot (SB)**: The programmer should compile the bootloader to support SB. On first boot, the ROM will generate a SB key to store in the eFuse. Then, the ROM generates an AES-based SHA digest over the bootloader using the SB key. The digest is stored in the flash. The programmer will also sign the firmware with an ECC private key, while the ECC public key is stored in the bootloader to verify the firmware.

A reliable and trusted secure facility can meet the provisioning requirements of the ESP32. Similar to the ECC608, **the ESP32 is fully capable of generating and storing its own private keys**, which significantly reduces the risk of exposure. As long as the build server has compiled the bootloader with the security features and the firmware has been signed, the secure facility only needs to upload these images to the flash, which will trigger the ESP32 to enable the security features. The secure facility can also perform some tests to check that security has been enabled. This ensures that the private keys are never exposed to anyone.

VI. EVALUATION

In this section, we discuss the area overhead of the ECC608 that is added to an MCU. We also explore the

Listing 4: Connecting to the DHT22.

```
#include "DHT22.h"
```

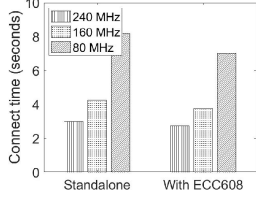



Fig. 7: AWS handshake time.

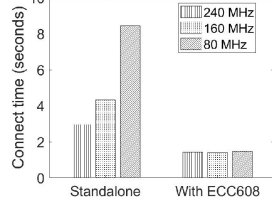


Fig. 8: EC2 handshake time.

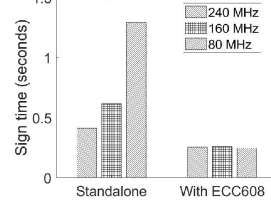


Fig. 9: ECDSA sig. gen. time.

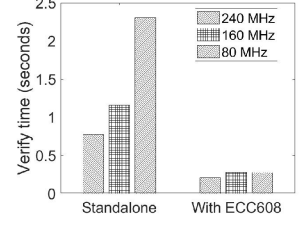


Fig. 10: ECDSA sig. verify time.

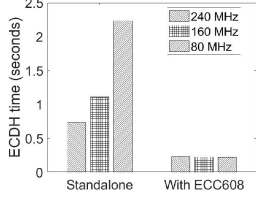


Fig. 11: ECDH time.

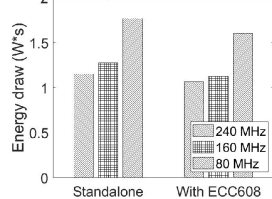


Fig. 12: AWS handshake energy draw.

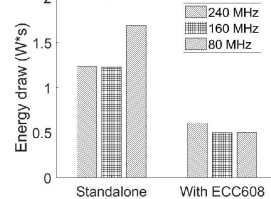


Fig. 13: EC2 handshake energy draw.

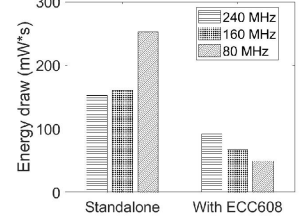


Fig. 14: ECDSA sig. gen. energy draw.

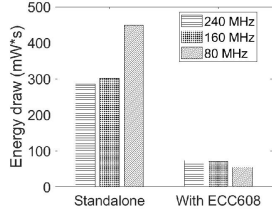


Fig. 15: ECDSA verify energy draw.

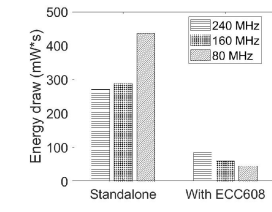


Fig. 16: ECDH energy draw.

improvements to the speed and energy consumption of the TLS handshake provided by the integrated ECC608 crypto chip. For performance assessment of the ESP32 security features, please refer to Appendix C.

A. ECC608 Area Overhead

PCB size is an important factor when considering IoT production costs. We have calculated the size of the ECC608 and WROOM and determined the area overhead of this crypto chip. The physical dimensions of the WROOM are roughly 459 mm^2 , while the ECC608 dimensions are about 29.4 mm^2 . This results in an area overhead of about 6.4% relative to the WROOM module. When considering the area of the overall circuit board, this shows that the area overhead of the ECC608 is quite minimal and will likely have an acceptable impact on production costs for IoT companies.

B. AWS Versus EC2

We have measured the network performance of *SIC*² on Amazon Web Services (AWS) IoT Core and Amazon Elastic Compute Cloud (EC2). **AWS IoT Core**, or simply AWS, is an IoT management cloud service. AWS can generate certificates for the end user that are signed by the Amazon Root CA. AWS also serves as an MQTT message broker, meaning end devices can connect to AWS using MQTT. This broker uses TLS on port 8883, allowing for a protected connection. Meanwhile, **EC2** is a service that allows users to configure and run virtual machines in the cloud. Our EC2 instance runs Ubuntu 18.04. To set up MQTT over TLS, we

used the Mosquitto software which can be used to establish an MQTT broker; Mosquitto can be configured to use TLS for mutual authentication and encryption, similar to AWS.

The difference in the network connection between AWS and EC2 lies in their server certificates. During the TLS handshake, AWS will present a server certificate signed by an RSA private key, while EC2 has been configured to use a certificate signed by an ECC private key. This means that during the TLS handshake, the ECC608 cannot be used to verify the AWS certificate and negotiate the session key, since ECC608 only supports ECC for public key cryptography. However, the ECC608 can be used to verify EC2's certificate and take advantage of the hardware acceleration.

C. ECC608 Speed

The ECC608 contains hardware acceleration of crypto operations, resulting in much better performance when compared to equivalent software implementations. We have measured the TLS handshake time between a remote server and a standalone ESP32 vs. one paired with the ECC608. We observe how clock speed impacts the handshake time by setting the ESP32 CPU speed to 240, 160, or 80 MHz. We also compare performance between AWS IoT and an EC2 server, the latter of which uses an ECC-based certificate and can perform ECDH with our ESP32. Each benchmark was executed 100 times, and we recorded the average runtime.

Figure 7 shows the total handshake time when connecting to AWS, while Figure 8 measures the EC2 handshake time. Connecting to AWS does not impact the connection time so drastically, since the ECC608 can only use the signature generation function to prove ownership of its certificate. However, when connecting to EC2, the handshake time reduces significantly, as much as 82% when the CPU clock speed is set to 80 MHz. This is because the ECC608 can also verify the server's certificate and perform ECDH. It can be observed that these operations form the majority of

computation, as the CPU clock speed has almost no impact on the handshake performance when the ECC608 is in use.

Figures 9 and 10 show metrics for ECDSA signature operations. In the worst case of 80 MHz, the ESP32 takes roughly 1.3 seconds to generate a signature and 2.3 seconds to verify a signature. By comparison, the ECC608 can consistently perform signature generation and verification in about 0.25 seconds.

Finally, we measure the time delay of ECDH which establishes the session key among the client and the server. Figure 11 shows these results. In the worst case of 80 MHz, the standalone ESP32 can perform ECDH in about 2.2 seconds. In comparison, the ECC608 reduces this latency to about 0.2 seconds. These results show that the hardware acceleration capabilities of the ECC608 can greatly benefit the networking performance of IoT applications.

D. Energy Consumption

To complement our performance metrics, we have also measured the energy consumption of ESP32 when performing the TLS handshake, ECDSA, and ECDH operations. Figures 12, 13, 14, 15, and 16 showcase these measurements. Note that the ECC608 itself also contributes to the total energy consumption, since it draws power from the ESP32. Despite this, our results indicate the ECC608 reduces energy consumption of the whole system. When using the crypto chip, the EC2 handshake requires only 0.6 Watt-seconds of energy, whereas a standalone ESP32 may require up to 1.7 Watt-seconds. At 80 MHz, the crypto chip reduces power usage by about 70%. ECDSA and ECDH benchmarks also reduce their individual energy consumption with the crypto chip. For instance, at 80 MHz, the ECC608 can perform the signature generation while drawing 49.8 megawatt-seconds, while the standalone ESP32 will draw 252 megawatt-seconds under this operation. These results are consistent with the signature verification and ECDH key exchange benchmarks.

VII. RELATED WORK

In this section, we discuss some software and hardware attacks that have targeted ESP32 and ESP8266 [27] in recent years.

Hardware Exploits. Researchers have exposed critical hardware vulnerabilities on ESP32-based smart devices. The LIFX Mini smart bulb was found to not implement flash encryption or secure boot, and JTAG was left completely open, leading to a full extraction of firmware details including WiFi credentials and a private RSA key [16]. A similar attack was performed on WIZ smart bulb [17]. Researchers also performed a voltage glitching attack on the ESP32 ROM with full security settings enabled, triggering a full readout of the security keys [15]. The latter attack cost several hundred dollars and could only be addressed with a hardware revision. [29]. By comparison, the crypto coprocessor discussed in this paper, the ATECC608A, employs several anti-tampering mechanisms to defend against such hardware exploits.

Software Exploits: Researchers have reported several vulnerabilities that affect ESP32 and ESP8266 software

libraries. The *Zero PMK Installation* vulnerability affects the EAP authentication framework [5]; attackers could force the Pairwise Master Key (PMK) to default to 0 and hijack a connection. In another vulnerability with the EAP framework, ESP32 will send an "EAPoL-Start" packet to the AP; if a malicious AP responds with a "success" packet, the ESP32 will crash [3]. In NONOS SDK (the official ESP8266 developer framework) 3.0 and earlier, the 802.11 MAC library fails to validate the bounds of the AuthKey Management (AKM) Suite Count value as well as the Pairwise Suite Count value. A malicious AP can send an arbitrarily large AKM packet and trigger a crash [4]. Note that ESP-IDF version 3.3 and NONOS version 3.1 address all of the aforementioned vulnerabilities. Carel Van Rooyen and Philipp Promeuschel [25] have shown that some ESP32 applications may be vulnerable to a stack-based buffer overflow attack if stack smashing protection is not enabled by the compiler. Our format string attacks differ from the above related works because our attacks are more general in nature (i.e., not specifically tied to any libraries), and the ESP32 provides no formal protection against format string attacks.

VIII. CONCLUSION

In this paper, we explore how cryptographic coprocessors may offer security protection to low-cost MCU based IoT devices by providing a hardware root of trust for private keys and a secure execution environment which is physically isolated from the host MCU. Software attacks are a major concern on IoT devices. We demonstrate two remote format string attacks on the popular ESP32 MCU. To thwart against these attacks, we pair the ESP32 with the ATECC608A crypto coprocessor, show how a manufacturing facility may provision private keys securely, and present implementation details on pairing the ESP32 with the ECC608. Finally, we show that the addition of a cryptographic coprocessor can advance the network performance of MCU based IoT devices by decreasing the TLS handshake time and energy consumption.

REFERENCES

- [1] M. Antonakakis, T. April, M. Bailey, M. Bernhard, E. Bursztein, J. Cochran, Z. Durumeric, J. A. Halderman, L. Invernizzi, M. Kallitsis, D. Kumar, C. Lever, Z. Ma, J. Mason, D. Menscher, C. Seaman, N. Sullivan, K. Thomas, and Y. Zhou. Understanding the mirai botnet. In *26th USENIX Security Symposium (USENIX Security 17)*, Aug. 2017.
- [2] C. Cimpanu. Hacker leaks passwords for more than 500,000 servers, routers, and iot devices. <https://www.zdnet.com/>, 2020.
- [3] M. E. Garbelini. Esp32/esp8266 eap client crash (cve-2019-12586). <https://matheus-garbelini.github.io/home/post/esp32-esp8266-eap-crash/>, 2019.
- [4] M. E. Garbelini. Esp8266 beacon frame crash (cve-2019-12588). <https://matheus-garbelini.github.io/home/post/esp8266-beacon-frame-crash/>, 2019.
- [5] M. E. Garbelini. Zero pmk installation (cve-2019-12587). <https://matheus-garbelini.github.io/home/post/zero-pmk-installation/>, 2019.
- [6] gosouth. Dht22 / am2302 library for esp32 (esp-idf). <https://github.com/gosouth/DHT22>, 2017.
- [7] A. Greenberg. The reaper iot botnet has already infected a million networks. <https://www.trendmicro.com/vinfo/pl/security/news/cybercrime-and-digital-threats/millions-of-networks-compromised-by-new-reaper-botnet>, October 2015.

- [8] A. Holdings. Arm mbed. <https://tls.mbed.org/>, 2020.
- [9] M. T. Inc. Atecc608a. <https://www.microchip.com/wwwproducts/en/ATECC608A>, 2018.
- [10] M. T. Inc. Atsamd21g18. <https://www.microchip.com/wwwproducts/en/ATsamd21g18>, 2020.
- [11] M. T. Inc. Cryptoauthlib - microchip cryptoauthentication library. <https://github.com/MicrochipTech/cryptoauthlib>, 2020.
- [12] T. Inc. Xtensa Instruction Set Architecture (ISA) Reference Manual.
- [13] A. Industries. Dht22 temperature-humidity sensor + extras. <https://www.adafruit.com/product/385>.
- [14] M. Korolov. What is a supply chain attack? why you should be wary of third-party providers. <https://www.csoonline.com/article/3191947/>, Jan 2019.
- [15] LimitedResults. Pwn the esp32 forever: Flash encryption and sec. boot keys extration. <https://limitedresults.com/2019/11/pwn-the-esp32-forever-flash-encryption-and-sec-boot-keys-extraction/>, 2019.
- [16] LimitedResults. Pwn the lifx mini white. <https://limitedresults.com/2019/01/pwn-the-lifx-mini-white/>, 2019.
- [17] LimitedResults. Pwn the wiz connected. <https://limitedresults.com/2019/02/pwn-the-wiz-connected/>, 2019.
- [18] R. LLC. Metasploit. <https://www.metasploit.com/>, 2020.
- [19] L. Luo, Y. Zhang, C. C. Zou, X. Shao, Z. Ling, and X. Fu. On runtime software security of trustzone-m based iot devices. <https://arxiv.org/abs/2007.05876>, Jul 2020.
- [20] Microsoft. What is azure sphere? <https://docs.microsoft.com/en-us/azure-sphere/product-overview/what-is-azure-sphere>, 03 2020.
- [21] C. Miller and C. Valasek. Remote exploitation of an unaltered passenger vehicle. <http://illmatics.com/Remote%20Car%20Hacking.pdf>, August 2015.
- [22] B. Pearson, C. Zou, Y. Zhang, Z. Ling, and X. Fu. *sic²*: Securing microcontroller based iot devices with low-cost crypto coprocessors. https://bpearson.net/papers/ICPADS_2020.pdf, 2020.
- [23] D. Rath. Open on-chip debugger. <http://openocd.org/files/thesis.pdf>.
- [24] R. Roemer, E. Buchanan, H. Shacham, and S. Savage. Return-oriented programming: Systems, languages, and applications. *ACM Trans. Inf. Syst. Secur.*, 15(1):2:1–2:34, 2012.
- [25] C. V. Rooyen and P. Promeuschel. Exploitation: Arm & xtensa compared. <http://nullcon.net/>, 2018.
- [26] E. Systems. Espressif achieves the 100-million target for iot chip shipments. https://www.espressif.com/en/news/Espressif_Achieves_the_Hundredmillion_Target_for_IoT_Chip_Shipments, Jan. 2018.
- [27] E. Systems. Esp8266 overview. <https://www.espressif.com/en/products/socs/esp8266/overview>, 2019.
- [28] E. Systems. Secure boot. <https://docs.espressif.com/projects/esp-idf/en/latest/security/secure-boot.html>, 2019.
- [29] E. Systems. Security advisory concerning fault injection and efuse protections (cve-2019-17391). https://www.espressif.com/en/news/Security_Advisory_Concerning_Fault_Injection_and_eFuse_Protections, 2019.
- [30] S. University. Format string vulnerability. http://www.cis.syr.edu/~wedu/Teaching/cis643/LectureNotes_New/Format_String.pdf.
- [31] ViperEye. Heap overflow: Vulnerability and heap internals explained. <https://resources.infosecinstitute.com/heap-overflow-vulnerability-and-heap-internals-explained/>, Jun. 2013.
- [32] B. Watters. Stack-based buffer overflow attacks: Explained and examples. <https://blog.rapid7.com/2019/02/19/stack-based-buffer-overflow-attacks-what-you-need-to-know/>, Feb. 2019.

APPENDIX A

FULL MCU DATASET

We conducted a survey on the state of the art for MCU hardware security, including the usage of crypto coprocessors, secure key storage, and trusted execution environments. Secure key storage means that the key is protected by a hardware root of trust, whereas trusted execution means that the crypto operations themselves are implemented in hardware and tamper-proof. Without such features, an attacker may confiscate secret data on an IoT device. Table I displays our full MCU dataset. In total, we surveyed 45

MCUs. Our results show that only 2 MCUs contain crypto coprocessors, the MT3620 by MediaTek and the CC3220S by TI. Additionally, only 9 MCUs offered secure storage of private keys, and only 28 offered any form of trusted execution. Of those MCUs which offer a trusted execution environment, only 6 offer it for ECDSA and ECDH crypto operations, only 3 offer it for RSA, and only 12 support any variant of SHA. These results show that many modern MCUs lack the capabilities of secure key storage and trusted execution environments that can provide hardware security.

APPENDIX B

ONE SCRIPT SECURITY FOR THE ESP32

If developers decide to utilize the ESP32 boards in their projects, they can take additional steps to protect hardware and firmware by incorporating flash encryption and secure boot into their applications. To complement this, we have provided a software solution for IoT device developers that may integrate the ESP32's security features into their projects. We label this software solution "One Script Security" or OSS. The flash encryption and secure boot functions of the ESP32 are enforced by the eFuse secure memory block as described in Section II.

Traditionally, in order to enable flash encryption and secure boot, developers can navigate through an interactive menu using the *idf.py* tool. This tool will configure various settings into a *sdkconfig* text file that is used to compile some features into the firmware. Moreover, developers can utilize the *espsecure.py* and *espefuse.py* to generate and set the flash encryption and secure boot keys. If the developer enables flash encryption or secure boot without first providing external keys, the ESP32 will internally generate these keys and they will be permanently inaccessible to the developer.

Internally generated security keys are considered the safest approach but can result in a soft brick of the device. For instance, the ESP32 supports only three reflashes of plaintext firmware while allowing for unlimited OTA updates; if the developer neglects to add OTA support into the firmware, then he cannot update the device any further. This restriction is enforced with the help of the *FLASH_ENCRYPT_CNT* eFuse. Without the secure boot key, developers cannot modify the software bootloader and produce a valid digest; however, many configuration options in *sdkconfig* may modify the bootloader without the user's knowledge, leading to accidental tampering of the chain of trust. Enabling flash encryption will also modify the bootloader, so these features must be enabled simultaneously.

We write OSS with the previous pitfalls in mind. **The software will externally generate the secure boot and flash encryption keys and read/write protect them to the current user.** The software will also encrypt and upload the firmware and bootloader at the proper address offsets. The secure boot key is derived from a digest of the private ECDSA *secure boot signing key*; in this way, the developer can always use the ECDSA private key to generate the secure boot key, which in turn will be used to generate a valid digest of the software bootloader. Since the flash encryption key

Device	Manufacturer	Processor	Crypto coprocessor	Secure storage	Trusted Execution	Other Security
A20	Marsboard	Cortex A7	No	No	No	No
A64	Allwinner	Cortex A53	No	Yes	Yes	Yes
AR9331	Atheros	MIPS 24K	No	No	No	No
BeagleBone Green	Seeed Studio	Cortex A8	No	No	No	No
BLE112	Silicon Labs	Intel 8081	No	No	No	No
CC2650	TI	Cortex M3	No	No	Yes	No
C3100M0D	TI	Cortex M3	No	No	Yes	Yes
CC3220S	TI	Cortex M4	Yes	Yes	Yes	Yes
CDXD5602	Sony	Cortex M4F	No	No	No	No
DM3725	TI	Cortex A8	No	No	No	No
eMote .Now	Samraksh	Cortex M3	No	No	No	No
Octa 5422	Exynos	Cortex A15	No	Yes	Yes	Yes
ATmega32U4	Atmel	AVR	No	No	No	Yes
H5	Allwinner	Cortex A53	No	No	Yes	No
i.MX 6SoloLite	NXP	Cortex A9	No	Yes	Yes	Yes
i.MX RT1060	NXP	Cortex M7	No	Yes	Yes	Yes
Jetson AGX Xavier	Nvidia	Arm V8	No	Yes	Yes	Yes
Jetson Nano	Nvidia	Arm A57	No	No	Yes	Yes
Kinetis KL8x	NXP	M0+	No	No	Yes	Yes
Kinetis MK20	NXP	Cortex M4	No	No	No	Yes
LPC5411	NXP	Cortex M4	No	No	No	Yes
MKW41Z	NXP	Cortex M0+	No	No	Yes	Yes
MSP430	TI	MSP430	No	No	Yes	No
MT3620	MediaTek	Cortex A7	Yes	Yes	Yes	Yes
MT7620n	MediaTek	MIPS 24KEc	No	No	No	No
MT7687F	MediaTek	Cortex M4	No	No	Yes	Yes
MT8163	MediaTek	Cortex A53	No	No	Yes	No
nRF52832	Nordic Semiconductor	Cortex M4	No	No	Yes	No
nRF52840	Nordic Semiconductor	Cortex M4	No	No	Yes	Yes
NuMicro M487	Nuvoton	Cortex M4F	No	No	Yes	Yes
Omega2S	Onion	MIPS 24K	No	No	No	No
Quark SE C1000	Intel	Quark	No	Yes	No	Yes
Quark SE D2000	Intel	Quark	No	Yes	No	Yes
RK3399	Rockship	Cortex A72	No	No	Yes	No
SAMD21	Microchip	Cortex M0+	No	No	No	Yes
SAMD5	Microchip	Cortex M4F	No	No	Yes	No
SAML11	Microchip	Cortex M23	No	No	Yes	Yes
SMART AT91RM9200	Microchip	ARM920T	No	No	No	No
STM32F0	ST	Cortex M0	No	No	No	Yes
STM32F2	ST	Cortex M3	No	No	Yes	No
STM32F7	ST	Cortex M7	No	No	Yes	Yes
STM32L5	ST	Cortex M33	No	No	Yes	Yes
STM32WB	ST	Cortex M4	No	Yes	Yes	Yes
X86 ULTRA	UDOO	Pentium N3710	No	No	Yes	Yes
Zynq-7000	Xilinx	Cortex A9	No	No	Yes	Yes
TOTAL	N/A	N/A	2 / 45 (4%)	9 / 45 (20%)	28 / 45 (62%)	27 / 45 (60%)

TABLE I: The full MCU dataset for hardware security of IoT devices.

is accessible to the developer, he can encrypt the firmware *before* uploading the ESP32; this can be done an unlimited number of times. As long as the flash encryption key and secure boot signing key are not deleted, a developer can always produce a valid software image. Additionally, OSS will disable the JTAG and UART debugging ports.

Importantly, our solution integrates components from ESP-IDF such as *idf.py* and *espsecure.py*, and depends on the general structure of the build directory (which stores plaintext elf and binary files). As a result, OSS can be imported into existing projects based on ESP-IDF, and developers do not have to migrate away from this platform. Furthermore, if the developer wishes, he may still upload firmware physically.

A. OSS Implementation

Our OSS solution is written as a shell script that can directly access components of ESP-IDF. The script assumes that ESP-IDF and the Xtensa toolchain have already been installed. The developer will effectively work in place of

idf.py flash, which is the typical method for building and uploading an app to the ESP32. The source code is available at Github ².

The functionality of OSS is split into two phases: *first time setup* and *repeat setup*. In *first time setup*, the ESP32 has not yet configured security settings and eFuses, while in *repeat setup*, the security settings have already been configured. OSS can predict which setup to execute by analyzing the ESP32's eFuse contents; if the flash encryption key and secure boot key read out "????...???", then this implies that the security keys have been set already and that the *repeat setup* is appropriate. On the author hand, a reading of "0000...0000" suggests that the security keys have not been set yet, and therefore *first time setup* should execute. Alternatively, we provide configuration options that allow the user to manually run the desired setup.

The *first time setup* works as follows. First, *espsecure.py*

²Available at <https://github.com/PBearson/ESP32-One-Script-Security>

will generate a new flash encryption key and secure boot signing key; these keys will be read/write protected to the current user. Then *espefuse.py* will burn and write-protect the following eFuses: *BLK1* (flash encryption key), *BLK2* (secure boot key), *FLASH_CRYPT_CNT*, *FLASH_CRYPT_CONFIG*, *ABS_DONE_0*, *JTAG_DISABLE*, *DISABLE_DL_ENCRYPT*, *DISABLE_DL_DECRYPT*, and *DISABLE_DL_CACHE*. Recall that the secure boot key is derived from the signing key. Finally, the script will proceed with building, signing, and encrypting the application before finally uploading it to the ESP32.

The *repeat setup* skips the process of configuring eFuses and proceeds directly to the build step. In this case, the bootloader is flashed to an offset of 0x0 rather than 0x1000, which accommodates for the newly added bootloader digest that is stored at the start of the bootloader image. To properly encrypt and upload applications, OSS must determine the appropriate address offset of the firmware, bootloader, partition table, and if applicable, the OTA data partition. The bootloader offset is strictly 0x0 or 0x1000, while the firmware and OTA data offsets can be determined from running *make partition_table*. The partition table offset can be determined from the *sdkconfig* text file.

To ensure that the bootloader and firmware will correctly compile in the security features into the binary, OSS will peak at the *sdkconfig*, and if the corresponding security settings have not been configured, OSS will not run. Without flash encryption support, an application will not have the necessary code to read or write to external flash over the SPI channel. Without secure boot support, the compiler will not generate the bootloader digest, and the bootloader will not contain the code necessary to validate the firmware.

APPENDIX C ESP32 SECURITY EVALUATION

We have evaluated the following characteristics relating to the security features of the ESP32: binary file size overhead; firmware build time; and run time.

We first evaluate the binary size overhead from incorporating flash encryption and secure boot into the application. The overhead arises from the firmware and bootloader needing to compile some additional functionalities into the image. To observe how the overhead may change with respect to app size, we have prepared a "large" app and a "small" app. The "small" app is a standard "hello world" program, while the "large" app implements a WiFi mesh/BLE client node.

Table II shows insight into the binary sizes when flash encryption and secure boot are enabled/disabled. In both the small app and large app, the bootloader binary increases only by 12 kB when security settings are enabled. In the small app, the firmware increases by 52.3 kB, while the large app increases by 34.7 kB. From these results, we can infer that the binary overhead of these security settings remains fairly constant with respect to the program size.

Next, we evaluate performance relating to the build time of the small app and large app. The build time encompasses the time required to compile, encrypt, and upload a program

TABLE II: Binary size insight (in *kilobytes*).

Binary	App size	Insecure	Secure
Bootloader	Small	25.38	37.39
Bootloader	Large	27.36	39.38
Firmware	Small	144.32	196.60
Firmware	Large	1407.10	1441.78

over UART. Note that when we run the full build process, *all* source files are compiled. In a real development scenario, most files do not need to re-compile or re-encrypt every time a developer updates the app.

TABLE III: Build time benchmarks (unit *seconds*).

Benchmark	App size	Insecure	Secure
Compile	Small	33.94	34.65
Compile	Large	52.90	53.74
Encrypt	Small	N/A	3.21
Encrypt	Large	N/A	19.53
Upload	Small	4.55	20.78
Upload	Large	25.03	134.87
All	Small	38.50	58.64
All	Large	77.92	208.14

Table III shows our evaluation metrics. It can be seen that in both the large app and small app, compilation time only increases by about one second when security is added to the chip. This correlates with the size difference of the image binaries. Encryption time and upload time increase linearly with respect to the application size. We found that the small application can be encrypted in about 3.2 seconds and uploaded in 20.8 seconds, while the large application takes 19.5 seconds to encrypt and 134.9 seconds to upload. The upload time for plaintext applications is considerably shorter. Our findings indicate that the full build process of a secure app can result in as much as a 167% time delay and up to 208 seconds or more, while smaller apps suffer from less delay. However, the resulting delay is acceptable, given that it will only occur at times when the developer must update the application.

TABLE IV: Run time benchmarks (unit μs).

Benchmark	Insecure	Secure
int8 add	146	147
int8 div	374	377
int32 add	109	109
int32 div	274	275
int64 add	445	447
int64 div	4478	4481
float sin	1.11×10^4	3.06×10^4
float div	2506	2510
double sin	1.10×10^4	$1.11e+4$
double div	3.16×10^4	3.16×10^4
double sqrt	4786	4793
copy string	3150	3150
matrix mult	1.57×10^5	1.57×10^5
read to flash	7.30×10^4	7.26×10^4
write to flash	5.23×10^7	4.77×10^7

Finally, we evaluate several run-time benchmarks on the ESP32 to measure the impact of flash encryption. The ESP32 contains internal flash encryption and decryption blocks which allow the internal SRAM to read and write to the encrypted flash memory over an SPI bus. As a result, we expect that the instruction throughput will incur some delay. To comprehensively assess the performance of the

ESP32, our benchmarks including addition and division of 8-bit integers, 32-bit integers, 64-bit integers, floating point numbers, and doubles, as well as the sine function using floats and doubles, the square root function, string copying, matrix multiplication, and finally, reading and writing to flash memory. The string copy benchmark would copy a string of length 512 bytes from one address to another, while the "read to flash" and "write to flash" benchmarks operated on 32 byte payloads. We ran each benchmark one thousand times at 240 MHz clock speed and recorded the total runtime to execute all instances of a benchmark.

Figure IV shows the results of each benchmark. It can be seen that the majority of workloads were not impacted by the flash encryption delay. We observed only a notable delay in the "float sin" benchmark, which increased by roughly 19 milliseconds. However, we also observe that the "write to flash" benchmark *decreased* by about 4.6 seconds. Our results indicate that the majority of computation-intensive workloads will not be impeded by flash encryption.

APPENDIX D MORE ATTACKS ON ESP32

A. Remote Cross Subroutine Attack

This attack exploits the Cross Subroutine Attack through the *sprintf()* vulnerability. Listing 5 provides pseudocode. The private key is defined in *handleReq()*. Upon receiving a GET request, *handleReq()* will call *sub1()*, which calls *sub2()*, which parses the GET request and performs the format string instruction before sending back a response.

Listing 5: Remote Cross Subroutine Attack pseudocode.

```
WebServer server(80);
void sub2() {
    char* res = malloc(64);
    char* param = server.arg(0).c_str();
    sprintf(res, param);
    server.send(200, "text/plain", res);
}
void sub1() { sub2(); }
void handleReq() {
    char key[32] = "THIS IS A PRIVATE KEY";
    sub1();
}
void setup() { server.on("/", handleReq); }
void loop() { server.handleClient(); }
```

This attack presents several new challenges. Namely, the ESP32 can support protection against the heap overflow attack. If the length of the format string exceeds the allocated length of the result buffer (64 bytes in this example), then the application will crash. Based on experiments, the stack pointer address of *sub2()* is 0x3ffb1de0, while the stack pointer of *handleReq()* is dumped from its register to address 0x3ffb1e04. To perform the typical attack, the user would pass the following GET request to the server:

`http://[IP addr]/?h=[Prefix]+%25s`

Here, *Prefix* indicates a sequence of "%x" characters. Unfortunately, when *sprintf()* formats this string, the length will exceed 64 bytes. This will trigger a system crash.

To overcome this challenge, an adversary can utilize the "n\$" format string character, which formats the n^{th}

parameter referenced by the format string. For example, the instruction `printf("%3$s", a, b, c)` will cast *c* as a string, while *a* and *b* will be skipped. As shown above, an adversary needs to access the contents of the fourteenth address pointed to by *sprintf()*. When utilizing the "n\$" character, the following GET request can be passed to the server:

`http://[IP addr]/?a=%2514%24s`

The formatted string will fit within the length of the result buffer. Additionally, we found that the private key is written 12 bytes after the *handleReq()* stack pointer. These bytes will be printed in addition to the private key, and a web browser may not be able to decode this response because some of the bytes are not valid ASCII characters. Instead, we use the *curl* command in Linux to send this HTTP request to the server and *hexdump* to view the raw HTTP response from the server.

B. Serial to TCP Attack

Even without the *sprintf()* vulnerability, a remote adversary can still observe the console output of an ESP32 via a **Serial to TCP** attack. In this attack, output from a serial device is redirected to an open TCP port. An adversary can connect to this TCP port to view all console output from the ESP32. This can lead to compromise of a private key.

The Serial to TCP attack occurs as follows. First, we assume that the ESP32 is vulnerable to the **printf()**-based Same Subroutine attack or Cross Subroutine Attack. We also assume that the ESP32 is connected to some victim server via USB so that its serial output can be read by the server. The victim server will write serial output from the ESP32 to some TCP connection. This can be done a variety of ways; in this example, we use the Linux *netcat* networking tool:

```
netcat -l [tcp_port] <
[serial_port] > [serial_port]
```

This command will forward console output from [serial_port] to the TCP connection at port [tcp_port]. An adversary can use tools like Metasploit [18] to discover vulnerabilities and run tools such as *netcat* on a remote server. When the TCP port is listening on the victim server, the remote adversary can connect to the port via *netcat* and trigger the format string attack as discussed earlier. The adversary will see *printf()* output dumped to their terminal and allow them to view the private key.