



**HARDWARE SECURITY
KNOWLEDGE AREA
(DRAFT FOR COMMENT)**

AUTHOR: Ingrid Verbauwhede – KU Leuven

EDITOR: Andrew Martin – Oxford University

George Danezis – University College London

REVIEWERS:

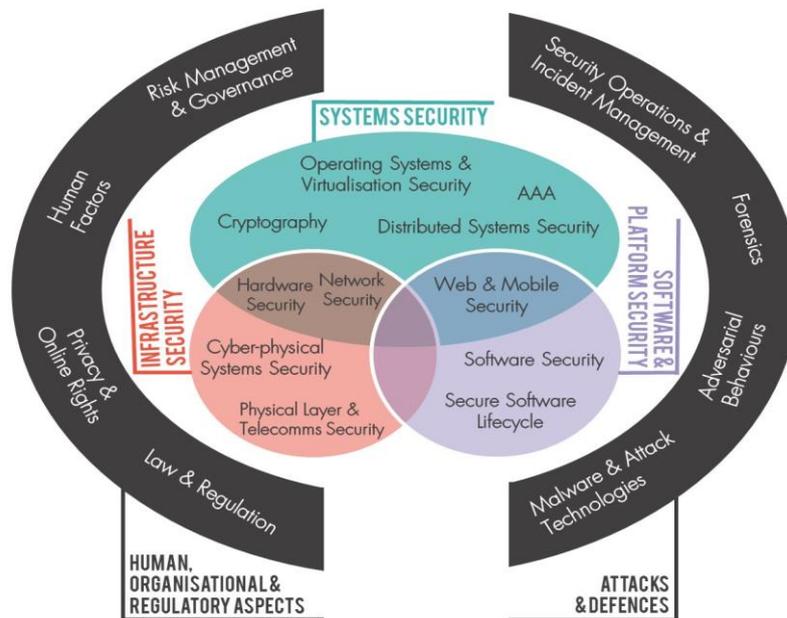
Srinivas Devadas – Massachusetts Institute of Technology

Paul England - Microsoft

Elisabeth Oswald – University of Bristol

Mark Ryan – University of Birmingham

Following wide community consultation with both academia and industry, 19 Knowledge Areas (KAs) have been identified to form the scope of the CyBOK (see diagram below). The Scope document provides an overview of these top-level KAs and the sub-topics that should be covered under each and can be found on the project website: <https://www.cybok.org/>.



We are seeking comments within the scope of the individual KA; readers should note that important related subjects such as risk or human factors have their own knowledge areas.

It should be noted that a fully-collated CyBOK document which includes issue 1.0 of all 19 Knowledge Areas is anticipated to be released in October 2019. This will likely include updated page layout and formatting of the individual Knowledge Areas.

Hardware Security

Ingrid Verbauwhede

August 2019

INTRODUCTION

Hardware security covers a broad range of topics from trusted computing to Trojan circuits. To classify these topics we follow the different hardware abstraction layers as introduced by the Y-chart of Gajski & Kuhn. These different layers of the hardware design process will be introduced in section 1, which explores the important concept of a root of trust and associated threat models. Next follows section 2 on measuring and evaluating hardware security. Subsequent sections gradually reduce the abstraction level. Section 3 describes secure platforms, i.e. a complete system or system-on-chip as trusted computing base. The discussion then moves into a section 4 covering hardware support for software security, namely, what features should a programmable processor include to support software security. This section is closely related to the Knowledge Area on software security. Register transfer level is the next abstraction level down, covered in section 5. Focus at this level is typically the efficient and secure implementation of cryptographic algorithms so that they can be mapped on ASIC or FPGA. This section is closely related to the Knowledge Area on cryptography. All implementations also need protection against physical attacks, most importantly against side-channel and fault attacks. Physical attacks and countermeasures are described in section 6. Section 7 describes entropy sources at the lowest abstraction level, close to CMOS technology. It includes the design of random numbers generators and physically unclonable functions. The last technical section describes aspects related to the hardware design process itself. This chapter ends with the conclusion and an outlook on hardware security.

1 Hardware design cycle and its link to hardware security

Hardware security is a very broad topic incorporating many other subject areas. In this section, these seemingly unrelated topics are grouped and ordered according to the design levels of abstraction as introduced by the Y-chart of Gajski & Kuhn [1]. While Gajski & Kuhn is a general approach to hardware design, in this chapter it is applied to the security aspects of hardware design and linked to threat models and associated root of trust.

Design abstraction layers are introduced in hardware design to reduce the complexity of the design. As indicated in 1, the lowest abstraction level a designer considers are the individual transistors at the center of the figure. These transistors are composed together to form basic logic gates, such as NAND, NOR gates or flip-flops, called the logic level. Going one abstraction layer up, at register transfer level, gates are grouped together to form modules, registers, ALU's, etc, and their operation is synchronised by a clock. These modules are then composed to form processors, specified by instruction sets, upon which applications and algorithms can be implemented.

By going up in the abstraction layers, details of underlying layers are hidden. This reduces design complexity at higher abstraction layers. The abstraction layers are represented by concentric circles in figure 1. Upon these circles, the Y-chart of Gajski & Kuhn introduces three design activities, represented by three axes: a behavioral axis, describing the behavior or *what* needs to be implemented (aka specifications); a structural axis describing *how* something is implemented and a

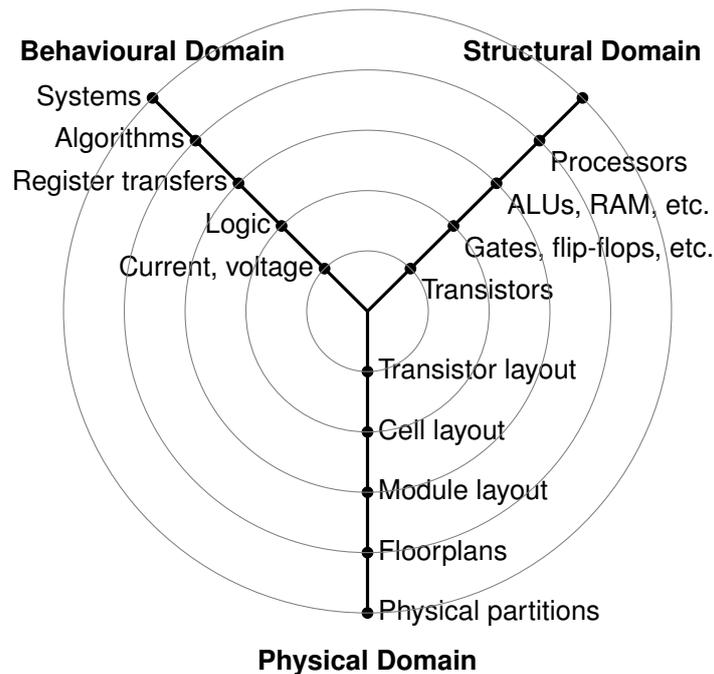


Figure 1: Gajski-Kuhn Y-chart

physical axis, pertaining to how the layouts are composed together at gate, module, chip and board level. An actual design activity is a ‘walk’ through this design space. Typically, one starts with the specifications at the top of the behavioral domain. These specifications (=what) are decomposed in components at the same level of abstraction (=how) moving from the behavioral axis to the structural axis. A structural component at one abstraction level becomes a behavioral component at one level down.

As an example of a walk through the design space: Assume a hardware designer is requested to implement a lightweight, low-power security protocol for an IOT (Internet of Things) device. This designer will only receive specifications on what needs to be designed: a security protocol aims to provide confidentiality and integrity (= what) and a set of cryptographic algorithms (= components) to support the protocol. The crypto-algorithms are provided as a behavioral specification to the hardware designer, who has the choice of implementing it as a dedicated co-processor, as an assembly program, or to support it with a set of custom instructions. Depending on costs and volumes, a choice of a target CMOS technology or an FPGA platform is made. This behavioral level will be translated into a more detailed register-transfer level description (e.g. VHDL or Verilog). At the RTL (Register Transfer Level), decisions need to be made as to whether this will be a parallel or sequential version, a dedicated or programmable design, with or without countermeasures against side-channel and fault attacks, etc.

Side channel attacks are attacks that are based on passive observations of a device operating under normal conditions. The attacker might monitor the execution time, the power consumption or electromagnetic radiation while the device is under normal operation. Since side-channel attacks are passive observations, the electronic circuit has very few or no means of realising that it is under attack. This is discussed in further detail in section 6.

1.1 Root of trust

A root of trust is a component used to realise a security function, upon which a designer relies but the trustworthiness of which cannot be explicitly verified [2]. It defines the trusted computing base and is defined by the trusted computing group as follows: “An entity can be trusted if it always behaves in

the expected manner for the intended purpose.” [3].

E.g. for a software or application developer, the root of trust is the TPM (Trusted Platform Module) or a SIM (subscriber identification module). For a hardware designer, the ultimate root of trust is the secure storage of the key in memory or the quality of the True Random Number Generator.

Hardware security is used as an enabler for software and system security. For this reason, hardware provides basic security services such as secure storage, isolation or attestation. The software or system considers the hardware as the trusted computing base. And thus from a systems or application view point, hardware has to behave as a trusted component. However, the hardware implementation can violate the trust assumption. E.g. Trojan circuits or side-channel attacks could leak the key or other sensitive data to an attacker. Hence, hardware itself also needs security. Moreover hardware needs security at all abstraction layers. Therefore, at every abstraction layer, a threat model and associated trust assumptions need to be made. An alternative definition for a root of trust in the context of design abstraction layers is:

“A root of trust is a component at a lower abstraction layer, upon which the system relies for its security.”

1.2 Threat model

A threat model is associated with each root of trust. When using a root of trust, it is assumed that threat model is not violated. This means that the threat model is also linked to the hardware abstraction layers. If we consider a root of trust at a particular abstraction layer, then all components that constitute this root of trust are also considered trusted.

Example 1: security protocols assume that the secret key is securely stored and not accessible to the attacker. The root of trust, upon which the protocol relies is the availability of secure memory to guard this key. For the protocol designer, this secure memory is a black box. The hardware designer has to decompose this requirement for a secure memory into a set of requirements at a lower abstraction layer. What type of memory will be used? On which busses will the key travel? Which other hardware components or software has access to the storage? Can there be side-channel leaks?

Example 2: It is during this translation of higher abstraction layer requirements from protocol or security application developers into lower abstraction layers for the hardware designers that many security vulnerabilities occur. Implementations of cryptographic algorithms used to be considered black boxes to the attacker: only inputs/outputs at the algorithm level are available to mount cryptanalysis attacks that are mostly mathematical. However, with the appearance of side-channel attacks this black box assumption no longer holds. Taking side-channel leakage into account, the attacker has the algorithm level information as well as the extra timing, power, and electro-magnetic information as observable from the outside of the chip. Thus the attacker model moves from black box to grey box. It is still assumed that the attacker does not know the details of the internals, e.g. the contents of the key registers.

Example 3: for programmable processors, the boundary between HW and SW is traditionally considered the instruction set architecture (ISA). The ISA is what is visible to the software programmer and the implementation of the ISA is left to the hardware designer. The ISA used to be considered the trust boundary for the SW designer. Yet, with the discovery of micro-architectural side-channel attacks, such as Spectre, Meltdown and Foreshadow, this ISA boundary is no longer black box, as micro-architectural information and leakages are available to the attacker [4].

1.3 Root of trust, threat model and hardware design abstraction layers

The decomposition in abstraction layers, in combination with electronic design automation (EDA) tools, is one of the main reasons that the exponential growth of Moore’s law was sustainable in the

past decades and it still is. This approach works well when optimising for performance, area, energy or power consumption. Yet for hardware security, no such general decomposition exists.

In this chapter, we propose to organise the different hardware security topics, their associated threat models and root of trust according to the hardware design abstraction layers, as there is no other known general body of knowledge available to organise the topics.

The advantage of this approach is that it can be used to identify the state of the art on different subtopics of hardware security. As an example, in the specific context of hardware implementations of cryptographic algorithms, the state of the art is well advanced and robust countermeasures exist to protect cryptographic implementations against a wide range of side-channel attacks, as shown in detail in section 5. Yet in the context of general processor security, e.g. to isolate process related data or to provide secure execution, new security hazards continue to be discovered on a regular basis.

<i>Abstraction level</i>	<i>Root of trust - functionality</i>	<i>Structural (how) - examples</i>	<i>Example Threats</i>	<i>Typical HW design activities</i>
System and application	Secure platforms	e.g. Trusted Execution (Trustzone, SGX, TEE), HSM, Secure Element	to support isolation, integrity, attestation, . . .	SW development
Processor	general purpose	e.g. shadow stack	SW vulnerabilities	HW/SW co-design
Processor	domain specific	Crypto specific RTL	Timing attacks	Behavioral synthesis
Register Transfer	Crypto specific	Building blocks,	SCA attack,	Logic synthesis
Logic	Resistance to SCA, Power, EM, fault	Masking, Circuit styles	SCA attack, fault	FPGA tools, standard cell design
Circuit and technology	Source of entropy	TRNG, PUF, Secure SRAM	Temp, glitches	SPICE simulations
Physical	Tamper Resistance	Shields, sensors	Probing, heating	Layout activities

Table 1: Design abstraction layers linked to threat models, root of trust and design activities

In an attempt to order the topics, table 1 summarises this organisation. The different abstraction layers are identified (first column) from a hardware perspective. The highest level (system and software) sits on top of the hardware platform. E.g. a system designer assumes that a secure platform is available. Thus the secure platform is the root of trust, providing security functionality. The second column describes the functionality provided by the root of trust. The third column describes how this functionality might be implemented. E.g. at the highest abstraction layer this might be by providing a Trusted Execution Module or a secure element. The fourth column describes the threat models and attack categories at that abstraction layer. For instance, at the system level, the system designer assumes to receive a module that provides isolation, integrity, attestation, etc. The last column describes typical design activities at this particular design abstraction layer.

This exercise is repeated for each abstraction layer and described in detail in each of the following sections.

At the processor level, one can distinguish general purpose programmable processors and domain

specific processors. General purpose processors should support a wide range of applications, which typically include software vulnerabilities. Hardware features are added to address these SW vulnerabilities, such as a shadow stack or measures to support hardware control flow integrity. Domain specific processors typically focus on a limited functionality. They are typically developed as co-processors in larger systems-on-chip. Typical examples are co-processors to support public key or secret key cryptographic algorithms. Time at the processor level is typically measured in instruction cycles.

Both, general purpose and domain specific processors, are composed together from computational units, multipliers and ALU's, memory and interconnect. These modules are typically described at the register transfer level: constant-time and resistance against side-channel attacks become the focus. Time at this level is typically measured in clock cycles.

Multipliers, ALU's, memories, interconnect and bus infrastructure are created from gates and flip-flops at the logic level. At this level, focus is on leakage through physical side-channels, power, EM, and fault attacks. Time is typically measured in absolute time (nsec) based on the available standard cell libraries or FPGA platforms.

The design of entropy sources requires knowledge and insights into the behaviour of transistors and the underlying CMOS technology. The design of these hardware security primitives is therefore positioned at the circuit and transistor level. Similarly, the design of sensors and shields against physical tampering require insight into the technology.

The layout of the table 1 does not aim to be complete. The idea is to illustrate each abstraction layer with an example. In the next sections, the hardware security goals and the associated threat model will be discussed in detail in relation to and relevant for each abstraction layer.

2 Measuring hardware security

Depending on the commercial application domain, several industrial and government organisations have issued standards or evaluation procedures. The most well known are the FIPS140-2, the common criteria CC evaluation and in the financial world the EMVCO. These evaluation procedures mostly focus on the implementation security of cryptographic algorithms. Note that this focus is narrower than the general hardware security requirement.

2.1 FIPS140-2

FIPS140-2 is a US NIST standard used for the evaluation of cryptographic modules. FIPS140-2 defines security levels from 1 to 4 (1 being the lowest). The following gives a description of the four levels from a physical hardware security point of view. Next to the physical requirements, there are also roles, services and authentication requirements (for more details see [5] and other Knowledge Areas).

Security level 1 only requires that an approved cryptographic algorithm be used, e.g. AES or SHA-3, but does not impose physical security requirements. Hence a software implementation could meet level 1. Level 2 requires a first level of tamper evidence. Level 3 also requires the tamper evidence, but additionally requires tamper resistance.

Tamper evidence means that there is a proof or testimony that tampering with a hardware module has happened. E.g. a broken seal indicates that a device was opened. A light sensor might observe that the lid of chip package was lifted.

Tamper resistance means that on top of tamper evidence, protection mechanisms are added to the device. E.g. if tampering is detected, then keys and sensitive material are reset.

Level 4 increases the requirements such that the cryptographic module can operate in physically unprotected environments. In this context, the physical side-channel attacks pose an important threat.

If any of these physical components depend on sensitive data being processed, information is leaked. Since the device is under normal operation, a classic tamper evidence mechanism will not realise that the device is under attack. See later in section 6.

2.2 Common criteria and EMVCo

"Common Criteria for information technology security evaluation" is an international standard for IT product security (ISO/IEC 15408), in short known as Common Criteria or CC. CC is a very generic procedure applicable to the security evaluation of IT products. Several parties are involved in this procedure. The customer will define a set of security specifications for its product. The manufacturer will design a product according to these specifications. An independent evaluation lab will verify if the product fulfills the claims made in the security requirements. Government certification bodies will issue a certification that the procedure was correctly followed and that evaluation lab indeed confirmed the claims made.

Depending on the amount of effort put into the security evaluation, the CC defines different Evaluation Assurance Levels (EAL). It ranges from basic functionally testing, corresponding to EAL1, to formally verified design and tested, corresponding to the highest level EAL7. CC further subdivides the process of evaluation into several classes, where most of the classes verify the conformity of the device under test. The 5th class (AVA) deals with the actual vulnerability assessment. It is the most important class as it searches for vulnerabilities and associated tests. It will assign a rating on the difficulty to execute the test, called the identification, and the possible benefit an attacker can gain from the penetration, called the exploitation. The difficulty is a function of the time required to perform the attack, the expertise of the attacker from layman to multiple experts, how much knowledge of the device is required from simple public information to detailed hardware source code, the number of samples required, and the cost and availability of equipment to perform the attack. A high difficulty level will result in a high score and a high level of the AVA class. The highest score one can obtain is an AVA level of 5, which is required to obtain a top EAL score.

Use of the Common Criteria process is well established in fields of smartcards and secure elements as they are used in telecom, financial, government ID's applications. It is also used in the field of Hardware Security Modules, Trusted Platform Modules and some more [6]. For certain classes of applications minimum sets of requirements are defined into protection profiles. Protection profiles are available for trusted platform modules (TPM), Javacards, Biometric passports, SIM cards, secure elements, etc.

Since certification comes from one government body, agreements exist between countries so that the certifications in one country are recognised in other countries. As an exception EMVCo is a private organisation that sets the specifications for worldwide interoperability of payment transactions. It has its own certification procedure similar to CC.

A good introduction to the topic can be found in [7] and a list of certified products on [6].

2.3 SESIP: Security Evaluation Standard for IoT Platforms

In the context of the Internet of Things (IOT) security evaluation, a recent initiative is the SESIP Security Evaluation scheme [8], currently at version 1.2. IOT devices are typically small, lightweight 'things' with limited accessibility via the internet. Several levels of threat model for IOT are possible: from only remote internet access, over various remote software attack options, to physical attack resistance. A comprehensive set of security functional requirements are defined: identification and attestation, product lifecycle, secure communication, software and physical attack resistance, cryptographic functionality including random number generation, and some compliance functionality to provide secure encrypted storage or provide reliable time. Similar to Common Criteria, SESIP provides several levels of assurance. Level 1 is the lowest level and consists of a self-assessment. The highest level of SESIP consists of a full CC evaluation similar to smart cards or secure elements. The levels in

between cover everything from a black box penetration testing over white box penetration, to testing with or without time limitations.

3 Secure Platforms

This section describes the goals and state-of-the-art features in secure platforms. At this high level of abstraction the system designer receives a complete chip or board as trusted computing base. The system designer assumes that the trusted root delivers a set of cryptographic functions, protected by the hardware and software inside the physical enclosure. Common to these platforms is that they are stand-alone pieces of silicon with a strict access policy. Depending on the provided functionality, the hardware tamper resistance and protection levels, and the communication interface, these secure platforms are used in different application fields (automotive, financial, telecom). The three most important platforms are the Hardware Security Module (HSM), the Subscriber Identification Module (SIM) and the Trusted Platform Module or TPM. These are briefly described next.

3.1 HSM Hardware Security Module

A HSM module will typically provide cryptographic operations, e.g. a set of public key and secret key algorithms, together with secure key management including secure generation, storage and deletion of keys. Essential to HSMs is that these operations occur in a hardened and tamper-resistant environment. A TRNG and a notion of a real-time clock is usually included. HSMs are mostly used in server back-end systems to manage keys or payment systems, e.g. in banking systems.

A HSM is used as a co-processor, attached to a host system. Its architecture typically includes a micro-processor/micro-controller, a set of crypto co-processors, secure volatile and non-volatile memory, TRNG, real-time clock, and I/O. The operations occur typically inside a tamper-resistant casing. In previous generations, multiple components reside on one board inside the casing. Recently, in some application domains, such as automotive, an HSM functionality is no longer provided as a stand-alone module but is now integrated as a secure co-processor in a larger System-on-Chip (SOC). Indeed Moore's law enables higher integration into one SOC.

For HSMs, compliance with the security levels is evaluated by specialised by independent evaluation labs.

3.2 Secure Element and Smartcard

Similar to an HSM, a Secure Element and a smart card provide a set of cryptographic algorithms, public key, secret key, HMAC, etc. together with secure key storage, generation and deletion. The main differences with HSMs are cost, size, and form factor. They are typically implemented as one single integrated circuit and have a much smaller form factor from around 50 cm² to less than 1 cm². The main difference between a smart card and a secure element sits in the form factor and the different markets they address. Secure elements are a more generic term, while smart cards have the very specific form factor of a banking card. They are produced in large volumes and need to be very cheap as they are used for SIM cards in cell phones and smart phones. They are also used in banking cards, pay-TV systems access cards, national identity cards and passports, and recently in IOT devices, vehicular systems and so on. Tamper resistance and physical protection are essential to secure elements. HSMs are a clear instance of what in computer architecture domain are called 'domain specific processors'.

A typical embedded secure element is one integrated circuit with no external components that consists of a small micro-controller with cryptographic co-processors, secure volatile and non-volatile storage, TRNG, etc. I/O is usually limited, through a specific set of pins, or through a NFC wire-less connection. Building a secure element is a challenge for a hardware designer, as one needs to combine security with non-security requirements of embedded circuits: small form factor (no exter-

nal memory), low power and/or low energy consumption in combination with tamper resistance and resistance against physical attacks, such as side-channel and fault attacks (see section 6).

3.3 Trusted Platform Module (TPM)

The TPM module has been defined by the Trusted Computing Group (TCG), an industry association, to provide specific security functions to the Personal Computer (PC) platform. More specifically, the TPM is a root of trust embedded on the PC platform, so that PC+TPM platform can identify itself and its current configuration and running software. [3]. The TPM provides three specific roots of trust: the root of trust for measurement (RTM), the root of trust for storage (RTS), the root of trust for reporting (RTR). Besides these three basic functions, other functionality of TPMs is being used: access to specific cryptographic functions, secure key storage, support for secure login, etc.

The TPM is implemented as a separate security module, much like a secure element but with a specific bus interface to a PC platform, e.g. through the LPC or I2C bus interface. Its architecture at minimum consists of an embedded micro-controller, several crypto coprocessors, secure volatile and non-volatile storage for root keys and a high quality true random number generator. It includes hardware engines for hash functions (SHA1 and SHA256), public key (RSA and ECC), secret key (AES) and HMAC calculations. Since a TPM is a separate module, physical protection and tamper resistance is essential for security.

The most recent TPM2.0 version broadens the application scope from PC oriented to also supporting networking, embedded, automotive, IOT, and so on. It also provides a more flexible approach in the functionality included. Four types of TPM are identified: the dedicated integrated circuit 'discrete element' TPM provides the highest security level. One step lower in protection level is the 'integrated TPM' as an IP module in a larger SOC. The lowest levels of protection are provided by the firmware and software TPM.

The adoption of TPM's has evolved differently from what was originally the focus of the TCG. Originally, the main focus was the support of a secure boot and the associated software stack, so that a complete measurement of the software installed could be made. The problem is that the complexity of this complete software base grows too quickly, making it too difficult to measure completely all the variations in valid configurations. Thus TPMs are less used to protecting a complete software stack up to the higher layers of software. Nonetheless, most new PCs now have TPMs but they are used to protect the encryption keys, avoid firmware roll-back, and assist the boot process in general.

4 Hardware support for SW security at architecture level

At the secure platform level, the complete module, i.e. HW and its enclosed embedded SW, are part of the trusted computing base. One level down on the abstraction layers, we make the assumption that all HW is trusted, while SW is no longer. Indeed, software vulnerabilities are a major source of security weaknesses (see the Knowledge Area on Software Security). To prevent the exploitation or to mitigate the effects of SW vulnerabilities, a large variety of HW modifications/additions to the processor architecture have been proposed in literature and have been included in commercial processors. We call this abstraction layer the HW/SW boundary: HW forms the trust boundary, while SW is no longer trusted. These security additions to the HW typically have a cost in extra area and loss in performance.

The most important security objectives at this level are to support protection, isolation and attestation for the software running on a processor platform [9], [10], [11].

- Protection: "A set of mechanisms for ensuring that multiple processes sharing the processor, memory, or I/O devices cannot interfere, intentionally or unintentionally, with one another by reading or writing each others' data. These mechanisms also isolate the operating system

from the user process" [9]. In a traditional computer architecture, usually the OS is part of the Trusted Computing Base (TCB), but the rest of the SW is not.

- With isolation, a hardware mechanism is added that controls access to pieces of software and associated data. Isolation separates two parties: a SW module might need protection from the surrounding SW, in which case, a Protected Module Architecture (PMA) provides a hardware guarantee that a piece of software runs unhindered from unwanted outside influences. In the opposite case, if we want to limit the effects of potentially tainted software to its environment, the software will be sandboxed or be placed into a 'compartment.' Protected Module Architectures are a hardware-only solution: the OS is not part of the TCB. More details are described in section 4.4
- With attestation, there is hardware support to demonstrate to a third party that the system, e.g. the code installed and/or running on a processor, is in a particular state. Attestation can be local or remote. Local attestation means that one SW module can attest its state to another one on the same computer platform. Remote attestation means that a third party, outside the computer platform can get some guarantee about the state of a processor.

In the context of general purpose computing, Virtual Machines (VMs) and Hypervisors have been introduced to support multiple operating systems on one physical processor. This sharing of resources improves efficiency and reuse. It can however only be realised by a safe and efficient sharing of physical memory: VMs should only be allowed to use the portions of physical memory assigned to them. Hence the important concept of *virtual memory* has been introduced: a technique that translates virtual addresses associated with one user to the physical address space which is shared between an Operating System and users. In this situation, different privilege levels are associated with different capabilities of the parties: the highest privilege level is associated with the hardware level. Supervisor, operating system and user applications have gradually lower privileges. More details can be found in the Knowledge Area on operating systems. The hardware supports protection by providing privileged instructions, control and status registers and sometimes support for multiple parallel threads.

In the context of embedded micro-controllers, with no operating system, and only one application, the hardware support could be limited to only machine level support. Memory protection could be added as an optional hardware module to the processor.

Other more advanced security objectives to support SW security might include:

- *Sealed storage* is the process of wrapping code and/or data with certain configuration, process or status values. Only under the correct configuration (e.g. program counter value, nonce, secret key, etc.) can the data be unsealed.
- *Dynamic root of trust* in combination with a *late launch* guarantees that even if the processor starts from an unknown state, it can enter a fixed known piece of code and known state. This typically requires special instructions to enter and exit the protected partition.
- *Memory protection* refers to the protection of data when it travels between the processor unit and the on-chip or off-chip memory. It protects against bus snooping or side-channel attacks or more active fault injection attacks.
- *Control flow integrity* is a security mechanism to prevent malware attacks from redirecting the flow of execution of a program. In hardware, the control flow of the program is compared on-the-fly at runtime with the expected control flow of the program.
- *Information flow analysis* is a security mechanism to follow the flow of sensitive data while it travels through the different processor components, from memory to cache over multiple

busses into register files and processing units and back. This is important in the context of micro-architectural and physical side-channel attacks.

In the next subsections a representative set of hardware approaches to address the above software security challenges are presented. Some hardware techniques address multiple security objectives. Some are large complex approaches, others are simple dedicated hardware features.

Note that a large body of knowledge on software-only approaches is available in literature. Mostly, they offer a weaker level of security as they are not rooted in a hardware (HW) root of trust. E.g. for control flow integrity, software only approaches might instrument the software code to check branches or jumps, while HW support might calculate MACs on the fly and compare these to stored associated MACs. Software (SW) only approaches are not further covered in this chapter but elsewhere (see the Knowledge Area on software security).

4.1 Trusted Execution Environment (TEE)

TEE was originally an initiative of Global Platform, a consortium of companies, designed to standardise a part of the processor as a trusted secure part. TEE has since evolved and in general, covers the hardware modifications made to processors to provide isolation and attestation to software applications. There is a large body of knowledge both from the industrial side as well as from the academic side.

TEE is a concept that provides a secure area of the main processor “to provide end-to-end security by protecting the execution of authenticated code, confidentiality, authenticity, privacy, system integrity and data access rights” [12]. It is important that the TEE is isolated from the so-called Rich Execution Environment (REE), which includes the untrusted OS. The reasoning behind this split is that it is impossible to guarantee secure execution and to avoid malware in the normal world due to the complexity of the OS and all other applications running there. The rich resources are accessible from the TEE, while the opposite is not possible. Global Platform does not specify the specifics on how these security properties should be implemented. Three main HW options are suggested. Option 1 assumes that every processor component on the IC can be split into a trusted and a rich part, i.e. the processor core, the crypto accelerators, the volatile and non-volatile memory are all split. Option 2 assumes that there is a separate secure co-processor area on the SOC with a well-defined HW interface to the rest of the SOC. Option 3 assumes a dedicated off-chip secure co-processor, much like a secure element.

Global Platform also defines a Common Criteria based protection profile (see section 3.2). It assumes the package of the integrated circuit is a black box [12] and thus secure storage is assumed by the fact that the secure asset remains inside the SOC. It follows the procedures of common criteria assurance package EAL 2 with some extra features. It pays extra attention to the evaluation of the random number generator and the concept of monotonic increasing time.

4.2 IBM 4758 Secure coprocessor

An early example, even before the appearance of the TEE of Global Platform is the IBM 4758 secure processor. Physical hardware security was essential for this processor: it contained a board with a general purpose processor, DRAM, separate battery backed-DRAM, Flash ROM, crypto accelerator (for DES), a random number generator and more. All of these components were enclosed in a box with tamper-resistant and tamper evidence measures. It was certified to FIPS 140-1, level 4 at that time [13].

4.3 ARM Trustzone

ARM Trustzone is a well known instantiation of a TEE. It is part of ARM processors integrated into system-on-chips (SOCs) mostly used for smartphones. The TEE is the secure part of the processor

and it runs a smaller trusted OS. It is isolated from the non-secure world, called the Rich Execution Environment, which runs the untrusted rich OS. The main HW feature to support this split is the Non-Secure (NS) bit. The AXI bus transactions are enhanced with an NS bit so that it can block the access of secure world resources by non-secure resources. Each AXI transaction comes with this bit set or reset. When the processor runs in the secure mode, the transaction comes with the NS bit set to zero, which gives it access to both secure and non-secure resources. When the processor runs in normal mode, it can only access resources from the normal world. This concept is extended to the level 1 and level 2 cache. These caches store an extra information bit to indicate if the code can be accessed by a secure or non-secure master. Special procedures are foreseen to jump from secure to non-secure and vice-versa. This is supported by a special monitor mode which exists in the secure world.

The split applied by ARM Trustzone is however a binary split. Applications from different vendors could co-exist together in the secure world, such that if one trusted component violates the system's security, the security can no longer be guaranteed. To address this issue, protected module architectures are introduced.

4.4 Protected Module Architectures and HWSW co-design solutions

If multiple software applications want to run on the same platform isolated from each other, then HW needs to isolate them from each other at a more fine granularity. This can be done by so-called protected module architectures. The basic idea is that small software modules can run protected from all other software running on the processor. And because they are small, their properties and behavior can be verified more thoroughly. The protection is provided by extra features added to the hardware in combination with an extremely small trusted software base if needed. In the Flicker project, the software TCB relies on only 250 lines of codes but requires a dedicated TPM chip [14]. Table 12 of the review work of [13], provides an in-depth comparison of several general purpose secure processor projects with their hardware and software TCB. The hardware TCB distinguishes between the complete mother board as TCB, e.g. for TPM usage, to CPU package only for SGX and other projects. The software TCB varies from a complete secure world as is the case for TrustZone to privileged containers in the case of SGX or a trusted hypervisor, OS or security monitor.

Even more advanced are solutions with a zero trusted software base: only the hardware is trusted. This is the case for the Sancus project [15]. It implements a program counter-based memory access control system. Extra hardware is provided to compare the current program counter with stored boundaries of the protected module. Access to data is only possible if the program counter is in the correct range of the code section. Progress of the program in the code section is also controlled by the hardware so that correct entry, progress and exit of the module can be guaranteed.

Intel's Software Guard Extensions (SGX) are also a protection mechanism at small granularity. Software modules of an application are placed in memory enclaves. Enclaves are defined in the address space of a process, but access to enclaves is restricted. Enclaves are created, initialised, and cleared by possibly untrusted system software, but operating in the enclave can only be done by the application software. Minimising the extra hardware to support SGX, and especially avoiding performance degradation is an important goal. The details of the hardware micro-architecture have not been disclosed: yet its most important parts are a memory encryption unit, a series of hardware enforced memory access checks and secure memory range registers [13].

4.5 Light-weight and individual solutions

The above listed solutions are mostly suited for general purpose computing, i.e. for platforms on which a complex SW stack will run. In literature, more solutions are proposed to provide extremely lightweight solutions to support specific security requests. SMART is one early example: it includes a small immutable piece of bootROM, considered the root of trust, to support remote attestation [16].

To protect against specific software attacks, more individual hardware countermeasures have been introduced. One example is a hardware shadow stack: to avoid buffer overflow attacks and to protect control flow integrity, return addresses are put on both the stack and the shadow stack. When a function loads a return address, the hardware will compare the return address of the stack to that of the shadow stack. They should agree for a correct return.

Another example is the protection of jump and return addresses to avoid buffer overflow attacks and other abuses of pointers. A simple but restrictive option is to use read-only memory, which fixes the pointer. A novel recent technique is the use of pointer authentication. The authentication code relies on cryptographic primitives. A challenge for these algorithms is that they should create the authentication tag with very low latency to fit into the critical path of a microprocessor. The ARMV8-A architectures therefore uses a dedicated low-latency crypto algorithm Qarma [17]. In this approach the unused bits in a 64-bit pointer are used to store a tag. This tag is calculated based on a key and on the program state, i.e. current address and function. These tags are calculated and verified on the fly.

Address Space Layout Randomisation is a general SW technique: its aim is to make it hard to predict the destination address of the jump. A detailed description can be found in the Knowledge Area on software security.

5 Hardware design for cryptographic algorithms at RTL level

The hardware features discussed so far are added to general purpose compute platforms, i.e. to a programmable micro-processor or micro-controller. General purpose means that a platform is created where the hardware designer does not know the future applications that will run on it. Flexibility, reflected in the instruction set, is then of importance. A second class of processors are domain-specific processors: they have limited or no programmability and are designed for one or a small class of applications.

When a dedicated processor is built for one or a class of cryptographic algorithms, this gives a lot of freedom to the hardware designer. Typically, the hardware designer will, starting from the cryptographic algorithm description, come up with hardware architectures at the Register Transfer Level (RTL) taking into account a set of constraints. Area is measured by gate count at RTL level. Throughput is measured by bits/sec. Power consumption is important for cooling purposes and measured in Watts. Energy, measured in Joules, is important for battery operated devices. It is often expressed in the amount of operations or amount of bits that can be processed per unit of energy. Hence the design goal is to maximise the operations/Joule or bits/Joule. The resistance to side channel attacks is measured by the number of measurements or samples required to disclose the key or other sensitive material. Flexibility and programmability are difficult to measure and are typically imposed by the application or class of applications that need to be supported: will the HW support only one or a few algorithms, encryption and/or decryption, modes of operation, initialisation, requirements for key storage, and so on.

A hardware architecture is typically described in a Hardware Description Language such as Verilog or VHDL. Starting from this description the two most important hardware platforms available to a HW designer are ASIC and FPGA. An Application Specific Integrated Circuit (ASIC) is a dedicated circuit fabricated in silicon. Once fabricated (baked) it cannot be modified anymore. FPGA or Field Programmable Gate Array is a special type of programmable device: it consists of regular arrays of 1-bit cells, that can be programmed by means of a bitstream. This special bitstream programs each cell to a specific function, e.g. a one bit addition, a register, a multiplexer, and so on. By changing the bit-stream the functionality of the FPGA changes. From the viewpoint of the register transfer level designer (RTL) the actual design process for either FPGA or ASIC doesn't differ that much. Similar design options are available: the designer can decide to go for serial or parallel architectures, making use of multiple design tricks to match the design with the requirements. The most well-known tricks

are to use pipelining to increase throughput, or unrolling to reduce latency, time multiplexing to reduce area, etc.

From implementation viewpoint, at this register transfer abstraction level, a large body of knowledge and a large set of Electronic Design Automation (EDA) tools exist to map an application onto a FPGA or ASIC platform [1]. Implementation results should be compared not only on the number of operations, but also on memory requirements (program memory and data memory), throughput and latency requirements, energy and power requirements, bandwidth requirements and the ease of which side-channel and fault attack countermeasures can be added. This large body of knowledge exists for implementations that focus on efficiency. However, when combining efficiency with security requirements, such as constant time execution or other countermeasures, there is a huge lack of supporting EDA tools (see section 8).

Cryptographic implementations are subdivided in several categories, enumerated below. The details of the cryptographic algorithms themselves are discussed in the Knowledge Area on cryptography. In this section, notes are made that are specific to the hardware implementations.

- Secret key algorithms: both block ciphers and stream ciphers result usually in compact and fast implementations. Feistel ciphers are chosen for very area constrained designs as the encryption and decryption hardware is the same. This is not the case for the AES algorithm for which encryption and decryption require different units.
- Secret key: lightweight algorithms. For embedded devices, over the years, many lightweight algorithms have been developed and implemented, e.g. Present, Prince, Rectangle, Simon or Speck cipher (see Knowledge Area on cryptography). Focus in these cases is mostly area cost. However, lightweight has lately been extended to include low power, low energy and especially low latency. Latency is defined as the time difference between input clear text and corresponding encrypted output or MAC. Having a short latency is important in real-time control systems, automotive, industrial IOT but also in memory encryption, and control flow integrity applications, etc. More knowledge will follow from the recent NIST call on light-weight crypto [18].
- Secret key: block ciphers by themselves are not directly applicable in security applications. They need to be combined with modes of operation to provide confidentiality or integrity, etc. (see Knowledge Area cryptography). In this context efficient implementations of authenticated encryption schemes are required: this is the topic of the CEASAR competition [19]. From an implementation viewpoint, the sequential nature of the authenticated encryption schemes makes it very difficult to obtain high throughputs as pipelining cannot directly be applied.
- Hash algorithms typically require a much larger area compared to secret key algorithms, especially the SHA3 algorithm, where different versions are large in area and slow in execution. Therefore, lightweight hash algorithms are a topic of active research.
- One important hardware application of hash functions is the mining of cryptocurrencies, such as Bitcoin, Ethereum, Litecoin and others, based on SHA2, SHA256, SHA3, etc. To obtain the required high throughputs, massive parallelism and pipelining is applied. This is however limited as hash algorithms are recursive algorithms. As a result, there is an upper bound on the amount of pipelining that can be applied [20]. Cryptocurrencies are discussed in more detail in the Knowledge Area on distributed systems.
- The computational complexity of public key algorithms is typically by 2 or 3 orders of magnitude higher than secret key and thus its implementation is 2 to 3 orders slower or larger. Especially for RSA and Elliptic curve implementations, a large body of knowledge is available, ranging from compact [21] to fast, for classic and newer curves [22].

- Algorithms resistant to attacks of quantum computers, aka post-quantum secure algorithms, are the next generation of algorithms requiring implementation in existing CMOS ASIC and FPGA technology. Computational bottlenecks are the large multiplier structures, with/without the Number Theoretic Transform, the large memory requirements and the requirements on random numbers that follow specific distributions. Currently, NIST is holding a competition on post-quantum cryptography [23]. Thus it is expected that after the algorithms are decided, implementations in hardware will follow.
- Currently, the most demanding implementations for cryptographic algorithms are those used in homomorphic encryption schemes: the computational complexity, the size of the multipliers and especially the large memory requirements are the challenges to address [24].

6 Side-channel attacks, fault attacks and countermeasures

This section first provides an overview of physical attacks on implementations of cryptographic algorithms. The second part discusses a wide range of countermeasures and some open research problems. Physical attacks, mostly side-channel and fault attacks, were originally of great concern to the developers of small devices that are in the hands of attackers, especially smart-cards and pay-TV systems. The importance of these attacks and countermeasures is growing as more electronic devices are easily accessible in the context of the Internet-of-Things (IOT).

6.1 Attacks

At the current state of knowledge, cryptographic algorithms have become very secure against mathematical and cryptanalytical attacks: this is certainly the case for algorithms that are standardised or that have received an extensive review in the open research literature. Currently, the weak link is mostly the implementation of algorithms in hardware and software. Information leaks also occur from the hardware implementation through side-channel and fault attacks. A distinction is made between passive or side-channel attacks versus active or fault attacks. A second distinction can be made based on the distance of the attacker to the device: attacks can occur remotely, close to the device, which is still non-invasive to actual invasive attacks. More details on several classes of attacks are below.

General side-channel attacks are passive observations of a compute platform. Through data dependent variations of execution time, power consumption or electromagnetic radiation of the device, the attacker can deduce information of secret internals. Variations of execution time, power consumption or electromagnetic radiations are typically picked up in close proximity of the device, while it is operated under normal conditions. It is important to note that the normal operation of the device is not disturbed. Thus the device is not aware that it is being attacked, which makes this attack quite powerful [25].

Side channel attacks based on variations on power consumption have been extensively studied. They are performed close to the device with access to the power supply or the power pins. One makes a distinction between Simple Power Analysis (SPA), differential and higher order power analysis (DPA), and template attacks. In SPA, the idea is to first study the target for features that depend on the key. E.g. a typical target in timing and power attacks are if-then-else branches that are dependent on key bits. In public key algorithm implementations, such as RSA or ECC, the algorithm runs sequentially through all key bits. When the if-branch takes more or less computation time than the else-branch this can be observed from outside the chip. SPA attacks are not limited to public key algorithms, they have also been applied to secret key algorithms, or algorithms to generate prime numbers (in case they need to remain secret). So with knowledge of the internal operation of the device, SPA is only required to collect one or a few traces for analysis.

With DPA, the attacker collects multiple traces, ranging from a few hundreds to millions in case of protected hardware implementations. In this situation, the attacker exploits the fact that the instant-

neous power consumption depends on the data that is processed. The same operation, depending on the same unknown sub-key, will result in different power consumption profiles if the data is different. The attacker will also build a statistical model of the device to estimate the power consumption as a function of the data and the different values of the subkey. Statistical analysis on these traces based on correlation analysis, mutual information and other statistical tests are applied to correlate the measured values to the statistical model.

Side channel attacks based on Electromagnetic radiations have been recognised early-on in the context of military communication and radio equipment. As a reaction, NATO and the governments of many countries have issued TEMPEST [26]. It consists of specifications on the protection of equipment against unintentional electromagnetic radiation but also against leakage of information through vibrations or sound. Electromagnetic radiation attacks can be mounted from a distance, as explained above, but also at close proximity to the integrated circuit. Electromagnetic probing on top of an integrated circuit can release very localised information of specific parts of an IC by using a 2D stepper and fine EM probers. Thus EM evaluation has the possibility to provide more fine grained leakage information compared to power measurements.

Timing attacks are another subclass of side-channel attacks [27]. When the execution time of a cryptographic calculation or a program handling sensitive data, varies as a function of the sensitive data, then this time difference can be picked up by the attacker. A timing attack can be as simple as a key dependent different execution time of an if-branch versus an else-branch in a finite state machine. Cache attacks, which abuse the time difference between a cache hit and a cache miss are an important class of timing attacks [28], [29], .

With a template attack, the attacker will first create a copy or template of the target device [30]. This template is used to study the behaviour of the device for all or a large set of inputs and secret data values. One or a few samples of the target device are then compared to the templates in the database to deduce secret information from the device. Template attacks are typically used when the original device has countermeasures against multiple executions. E.g. it might have an internal counter to log the number of failed attempts. Templates can be made based on timing, power or EM information. As machine learning and AI techniques become more powerful, so will the attack possibility with template attacks.

Processor architectures are very vulnerable to timing attacks. The problem of information leaks and the difficulty of confinement between programs was already identified early on in [31]. Later timing variations in cache hits and misses became an important class of timing attacks [32]. Recently gaining a lot of attention are the micro-architectural side-channel attacks, such as Spectre, Meltdown, Foreshadow. They are also based on the observation of timing differences [4][32]. The strength of the attacks sits in the fact that they can be mounted remotely from software. Modern processors include multiple optimisation techniques to boost performance not only with caches, but also speculative execution, out-of-order execution, branch predictors, etc. When multiple processes run on the same hardware platform, virtualisation and other SW techniques isolate the data of the different parties in separate memory locations. Yet, through the out-of-order execution or speculative execution (or many other variants) the hardware of the processor will access memory locations not intended for the process by means of so-called transient instructions. These instructions are executed but never committed. They have however touched memory locations, which might create side channel effects, such as variations in access time, and thus leak information.

Fault attacks are active manipulations of hardware compute platforms [33]. The result is that the computation itself or the program control flow is disturbed. Faulty or no outputs are released. Even if no output is released or the device resets itself, this decision might leak sensitive information. One famous example is published in [34]: it describes an RSA signature implementation which makes use of the Chinese Remainder Theorem (CRT). With one faulty and one correct result signature, and some simple mathematical calculations, the secret signing key can be derived. Physical fault-attacks could be a simple clock glitching, power glitching, heating up or cooling down a device. These require

close proximity to the device but are non-invasive.

With more expensive equipment, and with opening the lid of the integrated circuit or etching the silicon down, even more detailed information of the circuit can be obtained. Equipment that has been used include laser attacks, focused ion beam (FIB), a scanning electron microscope (SEM) and other high-end lab tools. This is typically equipment that has been designed for chip reliability and failure analysis.

6.2 Countermeasures

There are no generic countermeasures that resist all classes of side-channel attacks. Depending on the threat model (remote/local access, passive/active, etc.) and the assumptions made on the trusted computing base (i.e. what is and what is not included in the root of trust), countermeasures have been proposed at several levels of abstraction. The most important categories are summarised below.

To resist timing attacks, the first objective is to provide hardware that executes the application or program in constant time independent of secret inputs, keys and internal states. Depending on the time granularity of the measurement equipment of the attacker, constant time countermeasures also need to be more fine grained. At the processor architecture level, constant time means a constant number of instructions. At the RTL level, constant time means a constant number of clock cycles. At logic and circuit level, constant time means a constant logic depth or critical path independent of the input data. At instruction level, constant time can be obtained by balancing execution paths and adding dummy instructions. Sharing of resources, e.g. through caches, makes constant time implementations extremely difficult to obtain.

At RTL level, we need to make sure that all instructions run in the same number of clock cycles. dummy operations or dummy gates, depending on the granularity level. Providing constant time RTL level and gate level descriptions is however a challenge. For performance reasons, design tools, both hardware and software compilers, will synthesise away the dummy operations or logic which were added to balance the computations.

As many side-channel attacks rely on a large number of observations or samples, randomisation is a popular countermeasure. It is used to protect against power, EM and timing side-channel attacks. Randomisation is a technique that can be applied at algorithm level: it is especially popular for public key algorithms, which apply techniques such as scalar blinding, or message blinding [35]. Randomisation applied at register transfer and gate level is called masking. Masking schemes randomise intermediate values in the calculations so that their power consumption can no longer be linked with the internal secrets. A large set of papers on gate level masking schemes is available, ranging from simple Boolean masking to threshold implementations that are provable secure under certain leakage models [36]. Randomisation has been effective in practice especially with public key implementation protection measures. The protection of secret key algorithms by masking is more challenging. Some masking schemes require a huge amount of random numbers, others assume leakage models that don't always correspond to reality. In this context, novel cryptographic techniques summarised under the label leakage resilient cryptography, are developed that are inherently resistant against side-channel attacks [37, 38]. At this stage, there is still a gap between theory and practice. For more details on leakage resilient cryptography, please refer to the Knowledge Area on cryptography.

Hiding is another major class of countermeasures. The idea is to reduce the signal to noise ratio by reducing the signal strength. Shielding in the context of TEMPEST is one such example. Similarly, at gate level, reducing the power signature or EM signature of standard cells or logic modules, will increase the resistance against power or EM attacks. Simple techniques such as using a jittery or drifting clock, and large decoupling capacitances will also reduce the signal to noise ratio.

Sometimes solutions for leaking at one abstraction level, e.g. power side channels, can be addressed

at a different abstraction level. Therefore, if there is a risk that an encryption key leaks from an embedded device, a cryptographic protocol that changes the key at a sufficiently high frequency, will also avoid side-channel information leakage.

General purpose processors such as CPUs, GPUs, and micro-controllers can not be modified once fabricated. Protecting against micro-architectural attacks after fabrication by means of software patches and updates is extremely difficult and mostly at the cost of reduced performance [4]. Micro-code updates are also a form of software, i.e. firmware update and not a hardware update. The main difference is that the translation from instructions to micro-code is a company secret, so for the user it looks like a hardware update. Providing generic solutions to programmable hardware is a challenge as it is unknown beforehand which application will run. Solutions to this problem will be a combined effort between hardware and software techniques.

Protection against fault attacks is mostly based on redundancy either in space or in time and by adding checks based on coding, such as parity checks. The price is expensive as calculations are performed multiple times. One problem with adding redundancy is that it increases the attack surface of side-channels. Indeed, due to the redundant calculations, the attacker has more traces available to perform time, power or EM side-channel attacks [35].

Many types of circuit level sensors are added to integrated circuits. Examples are light sensors that detect that the lid of a package has been opened. Mesh metal sensors layered-out in top level metal layers can detect probing attacks. Temperature sensors detect heating or cooling of the integrated circuit. Antenna sensors to detect EM probes close to the surface have been developed: these sensors measure a change in EM fields. And sensors that detect manipulation of the power supply or clock can be added to the device. Note that adding sensors to detect active manipulation can again leak extra information to the side channel attacker.

Joint countermeasures against side-channel and fault attacks are challenging and an active area of research.

7 Entropy generating building blocks: random numbers, physically unclonable functions

Sources of entropy are essential for security and privacy protocols. In this section two important sources of entropy related to silicon technology are discussed: random number generators and physically unclonable functions.

7.1 Random number generation

Security and privacy rely on strong cryptographic algorithms and protocols. A source of entropy is essential in these protocols: random numbers are used to generate session keys, nonces, initialization vectors, to introduce freshness, etc. Random numbers are also used to create masks in masking countermeasures, random shares in multi party computation, zero-knowledge proofs, etc. In this section the focus is on cryptographically secure random numbers as used in security applications. Random numbers are also used outside cryptography, e.g. in gaming, lottery applications, stochastic simulations, etc.

An ideal RNG should generate all numbers with equal probability. Secondly, these numbers should be independent from previous or next numbers generated by the RNG, called forward and backward secrecy. The probability can be verified with statistical tests. Each standard includes a large set of statistical tests aimed at finding statistical weaknesses. Not being able to predict future values or derive previous values is important not only in many security applications, e.g. when this is used for key generation, but also in many gaming and lottery applications.

In general, random numbers are subdivided in two major classes: the pseudo-random number generator (PRNG) also called deterministic random bit generator (DRBG) and the true random number

generator (TRNG) or non-deterministic random bit generator (NRBG). The design, properties and testing of random numbers is described in detail by important standards issued in the US by NIST. NIST has issued the NIST800-90A for deterministic random number generators, the NIST800-90B for entropy sources, and NIST800-90C for random bit generation constructions [39], [40] [41]¹. In Germany and by extension in most of Europe, the German BSI has issued two important standards: the AIS 20 for functionality classes and evaluation criteria for deterministic random number generators and the AIS-31 for physical random number generators [42, 43, 44].

Pseudo-random number generators are deterministic algorithms that generate a sequence of bits or numbers that look random but are generated by a deterministic process. Since a PRNG is a deterministic process, when it starts with the same initial value, then the same sequence of numbers will be generated. Therefore it is essential that PRNG starts with a different start-up value each time the PRNG is initiated. This initial seed can either be generated by a slow true random number generated or at minimum by a non-repeating value, e.g. as provided by a monotonic increasing counter. A PRNG is called cryptographically secure if the attacker, who learns part of the sequence, is not able to compute any previous or future outputs. Cryptographically secure PRNGs rely on cryptographic algorithms to guarantee this forward and backward secrecy. Forward secrecy requires on top a regular reseeding to introduce new freshness into the generator. Hybrid RNG have an additional non-deterministic input to the PRNG.

PRNGs provide conditional security based on the computational complexity of the underlying cryptographic algorithms. See the Knowledge Area on cryptography for more details. In contrast, ideal true random number generators provide unconditional security as they are based on unpredictable physical phenomena. Thus their security is guaranteed independent of progress in mathematics and cryptanalysis.

The core of a true random number generator consists of an entropy source, which is a physical phenomena with a random behavior. In electronic circuits, noise or entropy sources are usually based on thermal noise, jitter and metastability. These noise sources are never perfect: the bits they generate might show bias or correlation or other variations. Hence they don't have full entropy. Therefore, they are typically followed by entropy extractors or conditioners. These building blocks improve the entropy per bit of output. But as the entropy extractor are deterministic processes, they cannot increase the total entropy. So the output length will be shorter than the input length.

Due to aging or environmental conditions, the quality of the generated numbers might vary or decrease over time. Therefore, the standards describe specific tests that should be applied at the start and continuously during the process of generating numbers. One can distinguish three main categories of tests. The first one is the total failure test, applied at the source of entropy. The second one are online health tests to monitor the quality of the entropy extractors. The third ones are tests for the post-processed bits. The requirements for these tests are well described in the different standards and specialised text books [45].

The challenge in designing TRNGs is first to provide a clear and convincing proof of the entropy source, second the design of online tests which at the same are compact and can detect a wide range of defects [46]. The topic of attacks, countermeasures and sensors for TRNGs, especially in the context of IOT and embedded devices, is an active research topic.

7.2 Physically Unclonable Functions

From a hardware perspective, Physically Unclonable Functions or PUFs, are circuits and techniques to derive unique features from silicon circuits, similar to human biometrics [47]. The manufacturing of silicon circuits results in unique process variations which cannot be physically cloned. The basic idea of PUFs is that these unique manufacturing features are magnified and digitised so that they can

¹NIST800-90C does not exist as a standard yet.

be used in security applications similar to the use of fingerprints or other biometrics. Process and physical variations such as doping fluctuations, line or edge widths of interconnect wires, result in variations of threshold voltages, transistor dimensions, capacitances, etc. Thus circuits are created that are sensitive to and amplify these variations.

The major security application for PUFs is to derive unique device specific keys, e.g. for usage in an IOT device or smart card. Traditionally, this storage of device unique keys is done in non-volatile memory, as the key has to remain in the chip even when the power is turned-off. However, non-volatile memory requires extra fabrication steps, which makes chips with non-volatile memory more expensive than regular standard CMOS chips. Thus PUFs are promised as a cheap alternative for secure non-volatile memory, because the unique silicon fingerprint is available without the extra processing steps. Indeed, each time the key is needed, it can be read from the post-processed PUF and directly used in security protocols. They can also replace fuses, which are large and their state is relatively easy to detect under a microscope.

The second security application is to use PUFs in identification applications, e.g. for access control or tracking of goods. The input to a PUF is called a challenge, the output the response. The ideal PUF has an exponential number of unique challenge response pairs, exponential in the number of circuit elements. The uniqueness of PUFs is measured by the inter-distance between different PUFs seeing the same challenge. The ideal PUF has stable responses: it replies with the same response, i.e. there is no noise in the responses. Moreover, PUF responses should be unpredictable and physically unclonable.

The ideal PUF unfortunately does not exist. In literature, two main classes of PUFs are defined, characterised by the number of challenge-response pairs they can generate. So-called weak PUFs are circuits with a finite number of elements, with each element providing a high amount of entropy. The number of possible challenge-response pairs grows typically linear with the area of the integrated circuit. Hence they are called weak PUFs. The most well-known example is the SRAM PUF [48]. These PUFs are typically used for key generation. The raw PUF output material is not directly usable for key generation as the PUF responses are affected by noise. Indeed, subsequent readings of the same PUF might result in slightly varying noisy responses, typically up to 20%. Thus after the entropy extraction follows secure sketch (similar to error correction) circuits to eliminate the noise and compress the entropy to generate a full entropy key [49]. The challenge for the PUF designer is to come up with process variations and circuits that can be used as key material, but which are not sensitive to transient noise. A second challenge is to keep all the post-processing modules compact so that the key-generation PUF can be included in embedded IOT devices.

The second class are the so-called strong PUFs. In this case, the number of challenge-response pairs grows large, ideally exponential, with the silicon area. The most well-known example is the arbiter PUF [50]. A small number of silicon elements are combined together, e.g. to create a chain of multiplexers or comparators, so that simple combinations of the elements create the large challenge-response space. Also in this case, the effects of noise in the circuits needs to be taken into account. Strong PUFs are promised to be useful in authentication applications, e.g. for access control. Each time a challenge is applied to the PUF, a response unique to the chip will be replied. The verifier will accept the response if it can be uniquely tied to the prover. This requires that the PUF responses are registered in a form of a database beforehand during an enrollment phase.

The problem with strong PUFs is that there is a strong correlation between different challenge-response pairs of most circuits proposed in literature. Hence all of these circuits are broken with machine learning techniques [51] and can not be used for authentication purposes. The fundamental problem is that very basic, mostly linear operations are used to combine PUF elements, which makes them easy targets for machine learning attacks. Ideally, these should be cryptographic or other computationally hard operations resistant to machine learning: unfortunately these cannot tolerate noise. Light-weight PUF based security protocols are an active area of research.

8 Hardware design process

In this section, several hardware security topics are described which are directly related to the lower design abstraction layers. One is the trust in the hardware design process itself. Directly related to this, is the problem of Trojan circuits. Also part of the hardware design process are circuit level techniques for camouflaging, logic locking, etc.

8.1 Design and fabrication of silicon integrated circuits

It is important to note that the hardware design process itself also needs to be trusted. Because of its design complexity, design at each abstraction layer relies on Electronic Design Automation (EDA) tools. The design, fabrication, packaging and test of silicon integrated circuits is an international engagement: silicon foundries are mostly located in Asia. Silicon design tools are most developed in the US, and silicon testing and packaging usually occur all over the world. For chips that end-up in critical infrastructure, such as telecommunication, military, aviation, trust and verification of the complete design cycle is essential.

Since silicon foundries and mask making are extremely expensive, very few countries and companies can still afford it and a huge consolidation is taking place in the industry. For critical infrastructure, governments demand more tools and techniques to increase the trustworthiness of this international design process. On this topic, large research projects are defined to come up with methods and tools to increase the trustworthiness of the design process and especially to assess the risk of Trojan insertions during the design process.

8.2 Trojan circuits

Trojan circuits are logic or gates added to large integrated circuits. As they are not part of the specified functionality, they are difficult to detect. They rely on the fact that they are extremely small in comparison with the large size of integrated circuits and SOCs. Trojan circuits are classified according to three main criteria [52, 53]. The first one is the physical characteristics of the Trojan, i.e. how is the Trojan inserted into the circuit. E.g. does it require logic modifications or only layout modifications. The second one is the activation characteristic: will the Trojan be turned on by an internal or external event, etc. The third characteristic classifies the type of action taken by the Trojan, e.g. will it leak information or will it destroy functionality, etc. The knowledge on this topic is summarised in [52, 53].

8.3 Circuit level techniques

To avoid visual inspection, circuit level *camouflaging* techniques are introduced [54]. These are standard cells or other modules that visually look the same, or they look camouflaged by random extra material. This is done to avoid visual inspection and reverse engineering based on visual inspection.

Another techniques to avoid loss of intellectual property is *logic locking* [55]. With this technique, extra gates are added to a circuit with a secret input. Only when the correct key is applied to the secret gates, will the circuit perform the correct functionality. This is an active research topic with logic locking schemes being proposed and attacked, with SAT solvers being a very useful tool in attacking the circuits.

8.4 Time

The concept of time and the concept of sequence of events are essential in security protocols. The TCG identifies three types of sequencing: a monotonic counter, a tick counter and actual trusted time [3]. A monotonic counter always increases, but the wall clock time between two increments is unknown. The tick counter increases with a set frequency. It only increases when the power is on. At power-off the tick counter will reset. Therefore the tick counter is linked with a nonce and methods are foreseen to link this with a real wall clock time. Trusted time is the most secure. It makes sure that

there is a link between the tick counter and the real wall clock time. From a hardware viewpoint it will require non-volatile memory, counters, crystals, continuous power, and an on chip clock generators. The connection to a real wall clock will require synchronisation and actual communication channel.

The importance of time is discussed in more detail in the Knowledge Area on security protocols.

9 Conclusion

Hardware security is a very broad topic, covering many different topics. In this chapter, a classification is made based on the different design abstraction layers. At each abstraction layer, the threat model, root of trust and security goals are identified.

Because of the growth of IOT, edge and cloud computing, the importance of hardware security is growing. Yet, in many cases hardware security is in conflict with other performance optimisations, such as low power or limited battery operated conditions. In these circumstances, performance optimisation is *the* most important design task. Yet it is also the most important cause of information leakage. This is the case at all abstraction layers: instruction level, architecture level and logic and circuit level.

Another trend is that hardware is becoming more 'soft'. This is an important trend in processor architecture, where FPGA functionality is added to processor architectures. The fundamental assumption that HW is immutable is lost here. This creates a whole new class of attacks.

A last big challenge for hardware security is the lack of EDA tools to support hardware security. EDA tools are made for performance optimization and security is usually an afterthought. An added challenge is that it is difficult to measure security and thus difficult to balance security versus area, throughput or power optimizations.

LINKS TO OTHER KA

The Knowledge Area on hardware security is linked with most other KA in cyber security. The most important ones are mentioned in the text.

- Knowledge Area Malware and attack technologies
- Knowledge Area Software security
- Knowledge Area Cryptography
- Knowledge Area in Distributed systems

KNOWN ABBREVIATIONS

ASIC = Application Specific Integrated Circuit

CC = Common Criteria

DRNG = Deterministic Random Number Generator (same as PRNG)

EAL = Evaluation Assurance Level

EDA = Electronic Design Automation

FPGA = Field Programmable Gate Array

HSM = Hardware Security Module

IC = Integrated Circuit

IOT = Internet of Things

KA = Knowledge Area
OS = Operating System
PC = Personal Computer
PRNG = Pseudo Random Number Generator (same as DRNG)
PUF = Physically Unclonable Function
ROT = Root of Trust
SOC = System on Chip
TCB = Trusted Computing Base
TEE = Trusted Execution Environment
TXT = Trusted Execution Technology
TCG = Trusted Computing Group
TPM = Trusted Platform Module
TRNG = True Random Number Generator
VLSI = Very Large Scale Integration
VM = Virtual Machine

REFERENCES

- [1] H. Kaislin, *Top-Down Digital VLSI Design: From Architectures to Gate-Level Circuits and FPGAs*. Morgan Kaufmann, 2015.
- [2] R. Anderson, *Security Engineering: A guide to Building Dependable Distributed Systems*. Wiley, 2001.
- [3] D. Grawrock, *Dynamics of a Trusted Platform: A building block approach*. Intel Press, 2008.
- [4] C. Canella, J. V. Bulck, M. Schwarz, M. Lipp, B. von Berg, P. Ortner, F. Piessens, D. Evtushkin, and D. Gruss, "A systematic evaluation of transient execution attacks and defenses," *CoRR*, vol. abs/1811.05441, 2018. [Online]. Available: <http://arxiv.org/abs/1811.05441>
- [5] National Institute of Standards and Technology, "FIPS 140-2 Security Requirements for Cryptographic Modules," may 25, 2001 (Change Notice 2, 12/3/2002).
- [6] List of Common Criteria products. [Online]. Available: <https://www.commoncriteriaportal.org/products/>
- [7] V. Lomne, "Common criteria certification of a smartcard:a technical overview," cHES 2016 Tutorial 1. [Online]. Available: <https://iacr.org/workshops/ches/ches2016/presentations/CHES16-Tutorial1.pdf>
- [8] "Security evaluation scheme for iot platforms, version 1.2," 2019. [Online]. Available: <https://www.trustcb.com/iot/sesip/>
- [9] D. A. Patterson and J. L. Hennessy, *Computer Organization and Design - The Hardware / Software Interface (5th Edition)*, ser. The Morgan Kaufmann Series in Computer Architecture and Design. Academic Press, 2014.
- [10] P. Maene, J. Goetzfried, R. D. Clercq, T. Mueller, F. Freiling, and I. Verbauwhede, "Hardware-Based Trusted Computing Architectures for Isolation and Attestation," *IEEE Transactions on Computers*, vol. 67, no. 3, pp. 361–374, 2017.
- [11] A. P. Martin, "The ten page introduction fo trusted computing," 2008.
- [12] Global Platform Device Committee, "EE Protection Profile," version 1.2, Public Release, November 2014, Document Reference: GPD_SPE_021.
- [13] V. Costan and S. Devadas, "Intel sgx explained," Cryptology ePrint Archive, Report 2016/086, 2016, <https://eprint.iacr.org/2016/086>.

- [14] J. M. McCune, B. J. Parno, A. Perrig, M. K. Reiter, and H. Isozaki, "Flicker: An execution infrastructure for tcb minimization," *SIGOPS Oper. Syst. Rev.*, vol. 42, no. 4, pp. 315–328, 2008.
- [15] J. Noorman, J. V. Bulck, J. T. Muehlberg, F. Piessens, P. Maene, B. Preneel, I. Verbauwhede, J. Goetzfried, T. Mueller, and F. Freiling, "Sancus 2.0: A Low-Cost Security Architecture for IoT Devices," *ACM Transactions on Privacy and Security*, vol. 20, no. 3, p. 33, 2017.
- [16] K. Eldefrawy, G. Tsudik, A. Francillon, and D. Perito, "SMART: secure and minimal architecture for (establishing dynamic) root of trust," in *19th Annual Network and Distributed System Security Symposium, NDSS 2012, San Diego, California, USA, February 5-8, 2012*, 2012.
- [17] R. Avanzi, "The qarma block cipher family – almost mds matrices over rings with zero divisors, nearly symmetric even-mansour constructions with non-involutory central rounds, and search heuristics for low-latency s-boxes," Cryptology ePrint Archive, Report 2016/444, 2016, <https://eprint.iacr.org/2016/444>.
- [18] NIST Lightweight Cryptography. [Online]. Available: <https://csrc.nist.gov/projects/lightweight-cryptography>
- [19] CAESAR: Competition for Authenticated Encryption: Security, Applicability, and Robustness. [Online]. Available: <https://competitions.cr.yp.to/caesar.html>
- [20] Y. K. Lee, H. Chan, and I. Verbauwhede, "Iteration Bound Analysis and Throughput Optimum Architecture of SHA-256 (384,512) for Hardware Implementations," vol. 4867, pp. 102–114, 2007.
- [21] Y. K. Lee, L. Batina, K. Sakiyama, and I. Verbauwhede, "Elliptic Curve Based Security Processor for RFID," *IEEE Transactions on Computers*, vol. 57, no. 11, pp. 1514–1527, 2008.
- [22] F. Turan and I. Verbauwhede, "Compact and Flexible FPGA Implementation of Ed25519 and X25519," *ACM Transactions on Embedded Computing Systems*, vol. 18, no. 3, p. 21, 2019.
- [23] NIST Post Quantum Cryptography. [Online]. Available: <https://csrc.nist.gov/Projects/Post-Quantum-Cryptography>
- [24] Homomorphic Encryption Standardization. [Online]. Available: <http://homomorphicencryption.org/introduction/>
- [25] S. Mangard, E. Oswald, and T. Popp, *Power Analysis Attacks: revealing the secrets of Smart Cards*. Springer-Verlag, 2007.
- [26] Link to NATO Tempest website. [Online]. Available: <https://www.ia.nato.int/niapc/tempest>
- [27] P. Kocher, "Timing attacks on implementations of diffie-hellman, rsa, dss, and other systems," in *Proceedings of the 16th Annual International Cryptology Conference on Advances in Cryptology*, ser. CRYPTO '96. Springer-Verlag, 1996, pp. 104–113.
- [28] D. Bernstein, "Cache-timing attacks on aes," 2005. [Online]. Available: <https://cr.yp.to/antiforgery/cachetiming-20050414.pdf>
- [29] D. A. Osvik, A. Shamir, and E. Tromer, "Cache attacks and countermeasures: The case of aes," in *Topics in Cryptology – CT-RSA 2006*. Springer Berlin Heidelberg, 2006, pp. 1–20.
- [30] S. Chari, J. Rao, and P. Rohatgi, "Template attacks," in *Cryptographic Hardware and Embedded Systems - CHES 2002*. Springer Berlin Heidelberg, 2003, pp. 13–28.
- [31] B. W. Lampson, "A note on the confinement problem," *Communications ACM*, vol. 16, no. 10, pp. 613–615, 1973.
- [32] D. Page, "MASCAB: a Micro-Architectural Side-Channel Attack Bibliography." [Online]. Available: <http://www.github.com/danpage/mascab>
- [33] D. Karaklajic, J.-M. Schmidt, and I. Verbauwhede, "Hardware Designer's Guide to Fault Attacks," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 21, no. 12, pp. 2295 – 2306, 2013.
- [34] D. Boneh, R. A. DeMillo, and R. J. Lipton, "On the importance of eliminating errors in cryptographic computations," *J. Cryptology*, vol. 14, no. 2, pp. 101–119, 2001.
- [35] J. Fan and I. Verbauwhede, "An Updated Survey on Secure ECC Implementations: Attacks, Countermeasures and Cost," in *Cryptography and Security: From Theory to Applications - Essays Dedicated to Jean-Jacques Quisquater on the Occasion of His 65th Birthday*, ser. Lecture Notes in Computer Science, D. Naccache, Ed., vol. 6805. Louvain-la-Neuve, Belgium: Springer-

- Verlag, 2012, pp. 265–282.
- [36] S. Nikova, C. Rechberger, and V. Rijmen, “Threshold Implementations Against Side-Channel Attacks and Glitches,” in *Information and Communications Security, 8th International Conference, ICICS 2006*, ser. Lecture Notes in Computer Science, N. Li, P. Ning, and S. Qing, Eds., vol. 4307. Raleigh, NC, USA: Springer-Verlag, 2006, pp. 529–545.
- [37] D. Bernstein, “Implementing “practical leakage-resilient symmetric cryptography”,” 2012, presented at CHES 2012 rumpsession. [Online]. Available: https://www.cosic.esat.kuleuven.be/ches2012/ches_rump/rs9.pdf
- [38] S. Belaïd, V. Grosso, and F. Standaert, “Masking and leakage-resilient primitives: One, the other(s) or both?” *Cryptography and Communications*, vol. 7, no. 1, pp. 163–184, 2015.
- [39] E. B. Barker and J. M. Kelsey, “Recommendation for random number generation using deterministic random bit generators, nist special publication 800-90a,” 2015. [Online]. Available: <https://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.800-90Ar1.pdf>
- [40] M. Turan, E. Barker, J. Kelsey, K. McKay, M. Baish, and M. Boyle, “Recommendation for the entropy sources used for random bit generation, nist special publication 800-90b,” 2018. [Online]. Available: <https://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.800-90B.pdf>
- [41] E. B. Barker and J. M. Kelsey, “Recommendation for random bit generator (rgb) constructions, second draft,” 2016. [Online]. Available: https://csrc.nist.gov/CSRC/media/Publications/sp/800-90c/draft/documents/sp800_90c_second_draft.pdf
- [42] “Functionality classes and evaluation methodology for deterministic random number generators, version 3,” 2013. [Online]. Available: https://www.bsi.bund.de/SharedDocs/Downloads/DE/BSI/Zertifizierung/Interpretationen/AIS_20_pdf.pdf
- [43] “Functionality classes and evaluation methodology for physical random number generators, version 3,” 2013. [Online]. Available: https://www.bsi.bund.de/SharedDocs/Downloads/DE/BSI/Zertifizierung/Interpretationen/AIS_31_pdf.pdf
- [44] W. Killmann and W. Schindler, 2011.
- [45] D. Johnston, *Random Number Generators - Principles and practices: A guide for engineers and programmers*. De Gruyter, 2018.
- [46] J. Balasch, F. Bernard, V. Fischer, M. Grujic, M. Laban, O. Petura, V. Rozic, G. V. Battum, I. Verbauwhede, M. Wakker, and B. Yang, “Design and Testing Methodologies for True Random Number Generators Towards Industry Certification,” in *International IEEE European Test Symposium - ETS 2018*, ser. IEEE Computer Society, Bremen, DE, 2018, p. 10.
- [47] R. Maes, *Physically Unclonable Functions: Constructions, Properties and Applications*, 1st ed. Springer Publishing Company, Incorporated, 2013.
- [48] J. Guajardo, S. S. Kumar, G.-J. Schrijen, and P. Tuyls, “Fpga intrinsic pufs and their use for ip protection,” in *Cryptographic Hardware and Embedded Systems - CHES 2007*. Springer Berlin Heidelberg, 2007, pp. 63–80.
- [49] Y. Dodis, R. Ostrovsky, L. Reyzin, and A. Smith, “Fuzzy extractors: How to generate strong keys from biometrics and other noisy data,” *SIAM J. Comput.*, vol. 38, no. 1, pp. 97–139, 2008.
- [50] B. Gassend, D. Clarke, M. van Dijk, and S. Devadas, “Silicon physical random functions,” in *Proceedings of the 9th ACM Conference on Computer and Communications Security*, ser. CCS ’02. ACM, 2002, pp. 148–160.
- [51] J. Delvaux, R. Peeters, D. Gu, and I. Verbauwhede, “A Survey on Lightweight Entity Authentication with Strong PUFs,” *ACM Computing Surveys*, vol. 48, no. 2, p. 42, 2015.
- [52] M. Tehranipoor and F. Koushanfar, “A survey of hardware trojan taxonomy and detection,” *IEEE Design Test of Computers*, vol. 27, no. 1, pp. 10–25, 2010.
- [53] R. Karri, J. Rajendran, K. Rosenfeld, and M. Tehranipoor, “Trustworthy hardware: Identifying and classifying hardware trojans,” *Computer*, vol. 43, no. 10, pp. 39–46, 2010.
- [54] J. Rajendran, M. Sam, O. Sinanoglu, and R. Karri, “Security analysis of integrated circuit camouflaging,” in *Proceedings of the 2013 ACM SIGSAC Conference on Computer & Communications Security*, ser. CCS ’13. ACM, 2013, pp. 709–720.

- [55] M. Yasin, A. Sengupta, M. T. Nabeel, M. Ashraf, J. J. Rajendran, and O. Sinanoglu, “Provably-secure logic locking: From theory to practice,” in *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, ser. CCS '17. ACM, 2017, pp. 1601–1618.