

Trusted Reference Monitors for Linux using Intel SGX Enclaves

Alexander Harri Bell-Thomas
Jesus College



*A dissertation submitted to the University of Cambridge
in partial fulfilment of the requirements for the degree of
Master of Engineering in Computer Science*

University of Cambridge
Computer Laboratory
William Gates Building
15 JJ Thomson Avenue
Cambridge CB3 0FD
UNITED KINGDOM

Email: Alexander.Bell-Thomas@cl.cam.ac.uk

June 6, 2020

Declaration

I, Alexander Harri Bell-Thomas of Jesus College, being a candidate the Part III of the Computer Science Tripos, hereby declare that this report and the work described in it are my own work, unaided except as may be specified below, and that the report does not contain material that has already been used to any substantial extent for a comparable purpose.

Total word count: 0

Signed:

Date:

This dissertation is copyright © 2020 Alexander Harri Bell-Thomas.
All trademarks used in this dissertation are hereby acknowledged.

Abstract

Write a summary of the whole thing. Make sure it fits in one page.

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 1 |
| 2 | Background | 5 |
| 2.1 | Information Flow Control | 5 |
| 2.1.1 | Motivation, History, and <i>Decentralised IFC</i> | 6 |
| 2.1.2 | Security Labels and Tags | 7 |
| 2.1.3 | Modelling | 7 |
| 2.2 | Intel SGX | 7 |
| 2.2.1 | Security Characteristics | 7 |
| 2.2.2 | Architecture and Implementation | 8 |
| 2.2.3 | Enclave Lifecycle | 10 |
| 2.2.4 | Attestation | 12 |
| 2.2.5 | Sealing | 15 |
| 2.2.6 | SGX Versions | 15 |
| 2.3 | Aspects of the <i>Linux</i> Kernel | 15 |
| 3 | Related Work | 17 |
| 4 | Design and Implementation | 19 |
| 5 | Evaluation | 21 |
| 6 | Summary and Conclusions | 23 |

List of Figures

| | | |
|-----|--|----|
| 2.1 | Abstract overview of SGX's protection in an adversarial environment. . | 8 |
| 2.2 | A high-level overview of the SGX hardware and software architecture. | 9 |
| 2.3 | The process of creating and initialising an enclave; details given in § 2.2.3. Purple components belong to the SGX platform. | 12 |
| 4.1 | Abstract overview of an SGX enclave's protections. | 20 |

List of Tables

| | | |
|-----|---|----|
| 2.1 | Overview of notable SGX x86 instructions in an enclave's lifecycle. | 13 |
|-----|---|----|

Chapter 1

Introduction

The task of defending computer systems against malicious programs and affording isolation to protected system components has always been exceedingly challenging to achieve. A system's *Trusted Computing Base*, or *TCB*, defines the minimal set of software, firmware, and hardware components critical to establish and maintain system security and integrity. This traditionally includes, amongst others; the OS kernel; device drivers; device firmware; and the hardware itself. Compromise of a trusted component inside a system's *TCB* is a direct threat to any secure application running on it. A common approach to hardening a system's security is to minimise its *TCB*, diminishing its potential *attack surface*.

A modern trend is to outsource the physical layer of a system to a foreign party, for example a *cloud provider* — this is beneficial both in terms of cost and flexibility, but confuses many security considerations which assume that the physical layer itself can be trusted. In this context there is no guarantee of this, as the physical layer is usually provided as a *virtual machine*, inflating the system's *TCB* with an external and transparent software layer, the underlying *hypervisor*.

The concept of *Trusted Execution Environments*, *TEEs*, has been explored by the security community for a very long time as potential protection against this, providing isolated processing contexts in which an operation can be securely executed irrespective of the rest of the system — one such example is software *enclaves*. *Enclaves* are

general-purpose *TEEs* provided by a CPU itself, protecting the logic found inside at the architectural level. Intel’s Software Guard Extensions (SGX) is the most prolific example of a *TEE*, affording a *black-box* environment and runtime for arbitrary apps to execute under. Introduced by Intel’s *Skylake* architecture, a partial view of the platform’s working can be found in whitepapers and previous publications.

SGX provides a *TEE* by enforcing the following:

1. Isolating a protected application, coined an enclave, from all other processes on the system at any privilege level using hardware enforcement.
2. Protecting reserved memory against attacks using a dedicated hardware component, the *Memory Encryption Engine (MEE)*.

... bit more here, talk about attestation and measurement briefly

A common usage pattern for enclaves in modern production systems is to build on top of a *libOS*.¹ This approach sees a trusted application built to depend on a modified operating system which is loaded alongside it in the enclave. Examples include *SGX-LKL*, *Graphene*, and *Occlum*. These projects allow SGX-unaware applications to be inexpensively ported into enclaves, but drastically inflate the *TCB* of the resulting program.

The aim of this work is to explore methods of hardening Linux with an SGX-driven *reference monitor* to track and protect host OS system resources using *information flow control* methods. Further, it aims to investigate whether bundling an entire operating environment into an enclave is a necessity, instead asking if simply using the host operating system, once hardened, would suffice in some situations.

Our Contribution

This work provides:

¹Library Operating System

- A prototype implementation of an enclave-based, modular *reference monitor*, empowering *information flow control* techniques to operate with autonomy and protection from the host operating system. Enforcement is achieved using a modified Linux kernel, with an overall *TCB* including only a minimal footprint of the core kernel alongside the enclave application.
- A userspace interposition library to near-transparently integrate unmodified applications to fully function under the new restrictions.
- A full port of the *libtomcrypt* cryptography library for use inside an SGX enclave.
- A rigorous investigation of the performance implications of this approach, featuring a lightly-modified version of the *Nginx* production webserver. Worst-case performance shows a 35% decrease in request throughput, with the common case reporting 7-11%. Additionally we report a median overhead of $39\mu s$ (IQR $26-72\mu s$, $n = 10^6$) per affected *system call*, matching or surpassing similar, non enclave-based, systems.

Chapter 2

Background

This chapter will cover a number of topics essential to understanding the rationale and implementation of the design as discussed in § 4. These include; an introduction to *Information Flow Control (IFC)*, the Intel SGX platform, and a overview of key aspects of the Linux kernel relevant to the architecture of the prototype.

2.1 Information Flow Control

IFC regulates how and where data is permitted to move and be transformed in a computer system. This differs from access control, which defines what resources may be used by an entity — IFC allows granular control over how they may be used once accessed, including restriction propagation between components.

Formally, IFC defines and enforces a non-interference policy between abstract *security contexts*. A simple, atomic example is the distinction between *unclassified* and *classified* data — here, information is only allowed to flow *up*, ensuring that an *unclassified* entity does not learn anything marked as *classified*. In general this form of relationship can be represented as a partial ordering over *security contexts*, enumerated in the form of a lattice.

Practical systems often require, however, dataflow adhering to a more complicated policy set — for example, supporting *declassification*. Work undertaken by Pasquier et al., [X] which will be the core influence of the IFC model developed in this project, constructs a pliable and efficient *decentralised information flow control (DIFC)* model suitable for provenance enforcement and auditing in the Linux kernel.

The notation and concepts adopted are primarily derived from Pasquier [X] and Krohn [Y]; any additional works will be individually cited but may use differing terminology.

2.1.1 Motivation, History, and *Decentralised IFC*

IFC has, in recent years, been increasing in support as a powerful methodology for ensuring granularly privacy whilst simultaneously not unduly restricting access to access to sensitive information. IFC annotates data records with opaque *labels*, referring to either their confidentiality or integrity status. Rather than simply restricting access to sensitive data, as would be the action taken by an access control mechanism, IFC opts to track data as it propagates — if an entity attempts to move this into an unknown, untrusted, or conflicting *security context* the IFC system prevents this to ensure data is not improperly released.

IFC originated from research conducted in the mid-1970's [X,Y] but has not, as of yet, seen mainstream adoption. A major reason for this is that early schemes were designed around the *multi-level security (MLS)* doctrine set out in the *Orange Book*: [Z] this locked IFC to a shallow set of broad labels, mirroring existing institutional segregation (such as *restricted*, *secret*, *top secret*). Policies were managed centrally, something easily applicable in settings with rigorous hierarchies such as the military, but unwieldy in organisation with diverse security procedures.

The majority of recent research in this areas has advocated *decentralised information flow control (DIFC)*, introduced by Myers and Liskov. [X] DIFC is more granular than schemes adhering to the *MLS* model; for example, creating two distinct *security contexts* for two files in the same folder. Policies are discretionary, allowing users to specify and modify the enforced policy for files they own.

2.1.2 Security Labels and Tags

2.1.3 Modelling

2.2 Intel SGX

Intel’s Software Guard Extensions, SGX, was first announced and detailed in a handful of whitepaper documents published in 2013. [X,Y,Z] It described a novel approach, creating in-CPU containers with their own protected memory pools. These regions, called *enclaves*, cannot be read from or written to by an unauthorised party due to fundamental protection mechanisms provided by the x86 architecture, even if running in *Ring 0*.¹ Figure 4.1 illustrates this. *Enclaves* guarantee both integrity and secrecy to the application running inside it, even in the presence of a malicious host.

Motivation At a high level SGX aims to achieve security for sensitive application by shielding them, and the resources it uses, against tampering and to provide a guarantee to end users about an enclave’s integrity; this is achieved using attestation and measurement (described in § 2.2.4). A driving use case is in a cloud computing context, where users are forced to trust a foreign party with both their data and business logic. By distributing encrypted, yet executable, containers targeting a single, unique SGX core, users can be assured that their information is safe, regardless of any virtualisation that may be taking place. Only the provisioned CPU is able to decrypt and execute the enclave, strictly in accordance with the restrictions of the SGX platform.

2.2.1 Security Characteristics

At its heart SGX is designed to be *trustworthy*; this is achieved in a number of ways, including robust enclaving provisioning, sealing and attestation. Intel enumerates SGX’s protections as follows;

- Memory security against observation and modification from outside the enclave; this is achieved using an in-die *Memory Encryption Engine (MEE)*, with a secret

¹x86 offers four protection *rings*, of which Linux uses two — *0* for the kernel, and *3* for userspace.

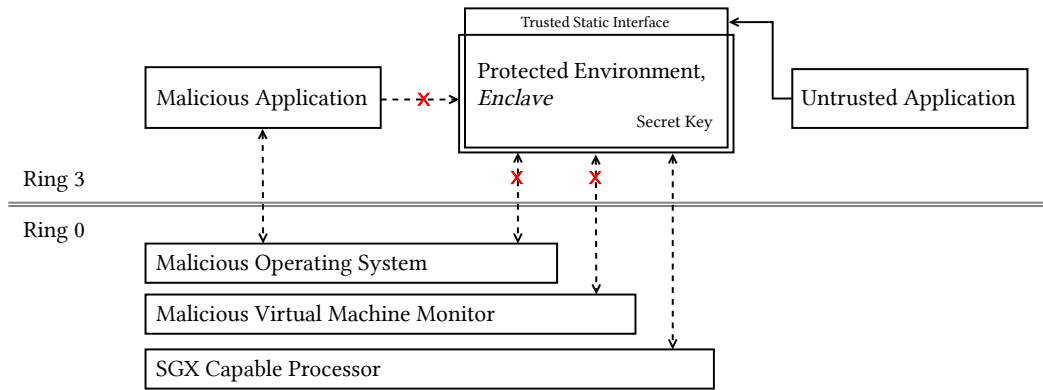


Figure 2.1: Abstract overview of SGX’s protection in an adversarial environment.

that rotates on every boot. This protection notably works against a host hypervisor, other enclaves, and anything running in supervisor mode.

- Attestation of an enclave to a challenger through the use of a permanent hardware security key for asymmetric encryption.
- Proxied software calls to prepare and transfer control in and out of an enclave. Arguments are securely marshalled according to a static enclave definition.
- SGX does not defend against reverse engineering or sidechannel attacks: [X,Y] this is the responsibility of the developer to mitigate.
- Debugging support is only provided via a specialised tool and only when an enclave is compiled with debugging enabled.

2.2.2 Architecture and Implementation

The SGX platform comprises a number of interlocking parts, as shown in Figure 2.2. Working from the hardware up, at the heart of the platform is the extended x86 instruction set and memory protection provided by an SGX-capable CPU.

Hardware Enclaves’ data and code is stored securely in *Processor Reserved Memory*, *PRM*; this is a set of pages in system memory that are presided over by the *MEE*. DMA² to *PRM* is always rejected. *PRM* consists of two datastructures; the *Enclave Page Cache*

²Direct Memory Access

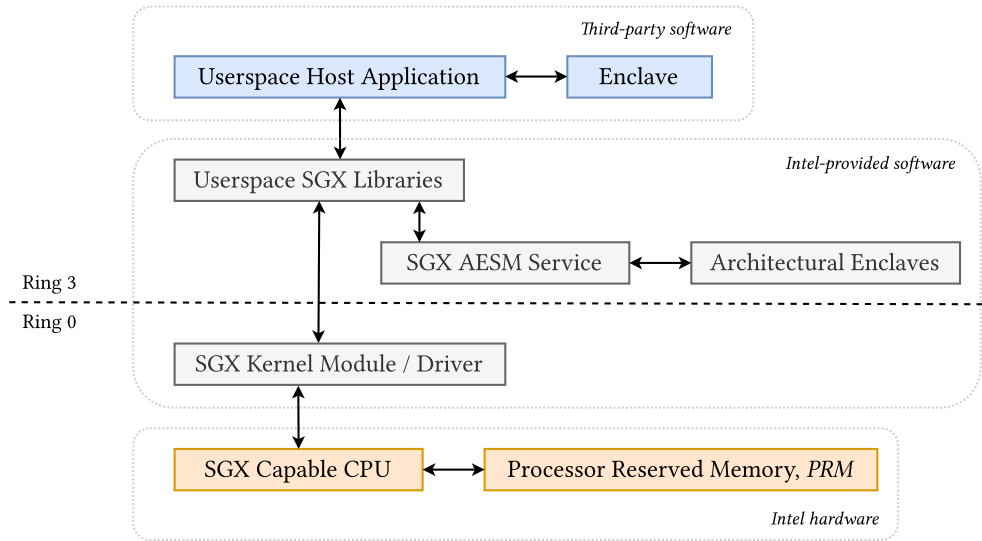


Figure 2.2: A high-level overview of the SGX hardware and software architecture.

Map (EPCM) and the *Enclave Page Cache (EPC)*. An individual enclave is defined by an *SGX Enclave Control Structure, SECS*; this is generated when an enclave is created and stored in a dedicated entry in the *EPC*. An enclave's *SECS* contains important information such as its (system) global identifier, its measurement hash and the amount of memory it is using. Access control information is stored in the *EPCM* alongside page validity flags, the owning enclave identifier and the page's type; this is not accessible from software. An attempt to resolve a page in *PRM* is successful only if the CPU is executing in enclave mode and its *EPCM* entry states it belongs the currently executing enclave – if this is not the case the lookup returns an unused page from generic system memory.

The host OS or hypervisor manages the *EPC* just as it does with normal system memory, swapping pages in and out according to its own policy, but must do so using SGX specific instructions. The *MEE* is responsible for ensuring the integrity and confidentiality of this process, encrypting and decrypting pages as they cross the *PRM* boundary. Data is verified with the use of an integrity tree, and encryption keys are generated at boot-time. Importantly the SGX architecture relies on the host OS being SGX-aware, empowering userspace applications to function without privilege; this is provided by the SGX driver, *isgx*.

Userspace services Starting an enclave requires retrieving a *launch token* from Intel’s *Launch Enclave*; this checks the signature and identity of the enclave to ensure it is valid. Access to the *Launch Enclave* and other architectural enclaves is provided by the AESM service; the userspace SGX libraries facilitate the communication mechanism. Other architectural enclaves include;

- The *Provisioning Enclave* — this verifies the authenticity of the platform and retrieves an enclave’s *attestation key* from the *Intel Provisioning Service*’s servers.
- The *Quoting Enclave* — this provides trust in the identity of the SGX environment and enclave being attested, by converting the locally generated *attestation key* to a remotely-verifiable *quote*.

Third-party enclaves Enclaves can only be entered via userspace, as detailed in § 2.2.3, and are always accompanied by a host application which acts as its untrusted counterpart. The host application calls the SGX SDK to build an enclave on its behalf using an enclave image, packaged as a standard shared library (`enclave.so`) and returns its *global identifier*. Control is passed from the host application to the enclave by invoking an enclave function via an *ECALL*. Execution flow can temporarily leave the enclave if it call one of the host application’s function via an *OCALL*. Execution naturally leaves enclave-mode when the *ECALL* terminates. Both *ECALLs* and *OCALLs* are defined statically in the enclave’s interface definition (`enclave.edl`). The necessary glue code is generated by the SGX SDK’s build toolchain at compile time; this ensures calls crossing the enclave boundary are marshalled safely and correctly.

2.2.3 Enclave Lifecycle

SGX instructions can be separated into two distinct groups; privileged and unprivileged. These, alongside a description of the function they perform, are enumerated in Table 2.1.³ The following description of the process of creating an enclave is illustrated in Figure 2.3.

³A handful of instructions not relevant to the explanation given here are omitted.

Preparing an enclave Execution begins with the host application; it needs to initiate the creation process, but must do so via a component with *Ring 0* privilege. This facility is provided by *isgx*, the SGX driver. The application first requests *isgx* to allocate the necessary number of pages to run the enclave $\langle 1 \rangle$;⁴ this is tracked and served from the driver's internal state $\langle 2 \rangle$.

The application continues by executing *ECREATE* with the metadata of the enclave to be loaded $\langle 3 \rangle$; the *MEE* checks that the pages being claimed are in fact vacant and populates the *SECS* page with the necessary information $\langle 4 \rangle$. Once this is complete the application prepares the remaining *EPC* pages using *EADD* $\langle 5 \rangle$ and loads the enclave's code and data $\langle 6 \rangle$.

At this point the enclave needs to be measured — the application calls *EEXTEND* $\langle 7 \rangle$, triggering the *MEE* to update the measurement hash in the *SECS* to aligns with the current state of the enclave's memory $\langle 8 \rangle$. Once the *EPC* memory is prepared the application requests for it to be finalised using *EINIT* $\langle 9 \rangle$: this operation requires the application to retrieve the *EINITTOKEN* from the *Launch Enclave*, locking the execution of the measured enclave to the CPU the token is generated on. Notably, pages cannot be added after *EINIT*,⁵ and an enclave cannot be attested to or entered before it. Finally, the initialised flag is set in the *SECS* and the enclave's hash updated for the final time $\langle 10 \rangle$.

Stepping into the enclave Once an enclave is created it can be invoked using the *EENTER* instruction. This can only jump to code explicitly defined in the enclave's interface definition and switches the CPU core to enclave mode. SGX uses a flag in the CPU core's *Thread Control Block* to prevent any other logical core following the current one into the enclave.

Interrupts and exceptions can be served to the enclave, just as with any other application. When in enclave mode control is not immediately passed over to the defined handler, but instead the enclave's current state is saved and cleared to ensure no data is

⁴These numbers correspond to events in Figure 2.3.

⁵This is only strictly true in SGXv1, as explained in § 2.2.6.

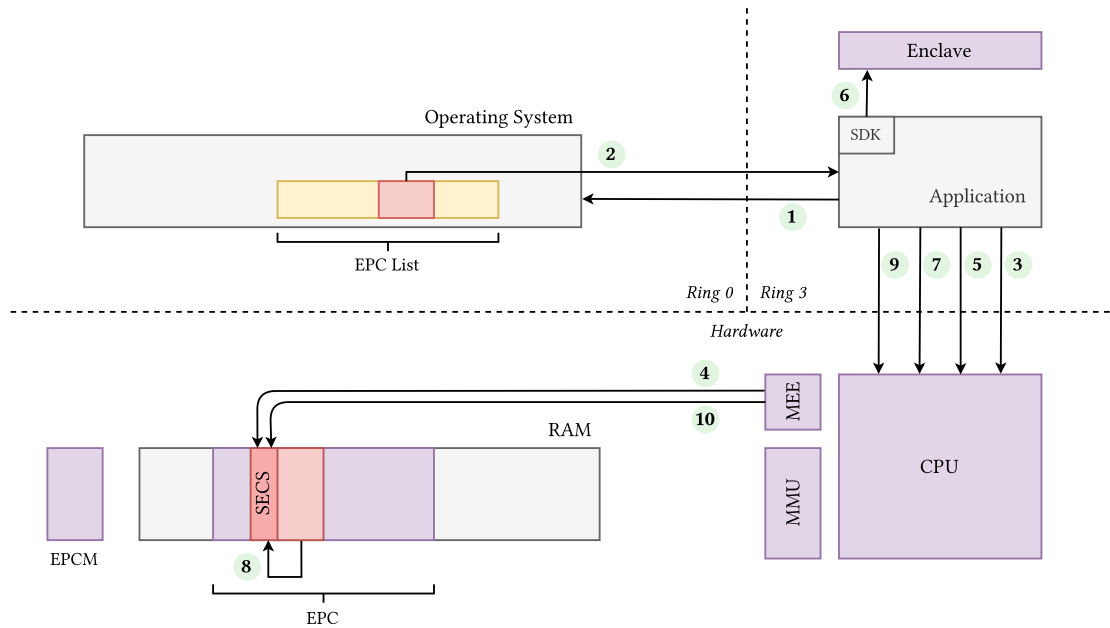


Figure 2.3: The process of creating and initialising an enclave; details given in § 2.2.3. Purple components belong to the SGX platform.

leaked. The *Asynchronous Enclave Exit* routine is then invoked and enclave mode disabled. Execution post-interruption is restarted with the `ERESUME` instruction. Once an enclave has finished executing the registers are erased and `EEXIT` called. Enclaves are terminated using the `EREMOVE` command; all claimed *EPC* pages are marked as invalid and the *SECS* page deleted.

A significant design decision made in the SGX architecture is that enclaves cannot be entered by a process operating in *Ring 0*; the required instructions simply aren't available. This forces all host applications to run in userspace, making interoperation with the kernel challenging, as will be discussed in § 4.

2.2.4 Attestation

An essential feature of the *trusted computing* model SGX creates is attestation, the process of verifying both the authenticity and integrity of component cryptographically. SGX achieves by creating two hashed values, or *signing identifies*, per enclave; `MRENCLAVE` and `MRSIGNER`.

| Execution Mode | Instruction | Function |
|----------------|-------------|---|
| Ring 0 | ECREATE | Generate and copy the <i>SECS</i> structure to a new page in the <i>EPC</i> , initialising a new enclave. |
| | EADD | Add a new <i>EPC</i> page for the current enclave; this is used to load initial code and data. |
| | EEXTEND | Updates the enclave's measurement during attestation; modifies the <i>SECS</i> . |
| | EINIT | The terminal instruction in an enclave's initialisation, finalising its attributes and measurement. |
| | EREMOVE | Permanently remove a page from the <i>EPC</i> ; usually invoked during enclave destruction. |
| Ring 3 | EENTER | Transfer control from the host application to a pre-determined location in an enclave. |
| | ERESUME | Re-enter the enclave after an interrupt/exception and resume execution. |
| | EEXIT | Restore the original operating mode at the location EENTER was triggered and flush the TLB. |
| | EGETKEY | Access platform cryptography keys required for attestation and sealing. |
| | EREPORT | Generate a <i>report</i> for an enclave's <i>attestation key</i> for an attestation process. |

Table 2.1: Overview of notable SGX x86 instructions in an enclave's lifecycle.

MRENCLAVE acts as a unique identifier for the contents of an enclave. It is generated by hashing the instructions and data passed when creating the enclave with ECREATE, EADD, and EEXTEND; the value is finalised and stored in the SECS on EINIT. This value depends on the exact content and ordering of the enclave's EPC pages. As long as the enclave's source remains the same, so will its MRENCLAVE.

MRSIGNER, also known as the enclave's *Sealing Identity*, is generated during the enclave build process — all production enclaves need to be signed using an RSA key provided by the compiling user (the *Sealing Authority*). The public key from this pair is stored in SIGSTRUCT, the *Enclave Signature Structure*. During an enclave's launch its SIGSTRUCT, which holds the signed compile-time MRENCLAVE value, is decrypted and crossreferenced with a freshly-computed runtime MRENCLAVE value to detect tampering. MRSIGNER is the same for all enclaves signed by the same *Sealing Authority*.

Local Attestation Two enclaves resident on the same system are able to attest their identities to each other using their MRENCLAVE and MRSIGNER values; this usually precedes the establishment of a shared secret (using a variant of *Diffie-Hellman* backed by the platform's master SGX key)⁶ for confidential communication between them.

Remote Attestation In addition to attestation between entities on the same platform, the Intel specification also provides a workflow for an enclave to attest its identity to a remote party. The system's *Quoting Enclave* verifies an enclave's local *quote* and creates a digital signature of it using the CPU's permanent hardware SGX private key. Through the use of an *Intel Enhanced Privacy Identifier (EPID)* this process can be carried out anonymously; it relies on information encoded in the CPU during the manufacturing process. The *Provisioning Enclave* assists in this process, especially as production enclaves are required to attest with Intel's provisioning service before executing. Remote attestation is not explicitly required in this project's architecture hence will not be covered in any further detail.

⁶Note for Harri: must check these details.

2.2.5 Sealing

2.2.6 SGX Versions

2.3 Aspects of the *Linux* Kernel

Chapter 3

Related Work

Chapter 4

Design and Implementation

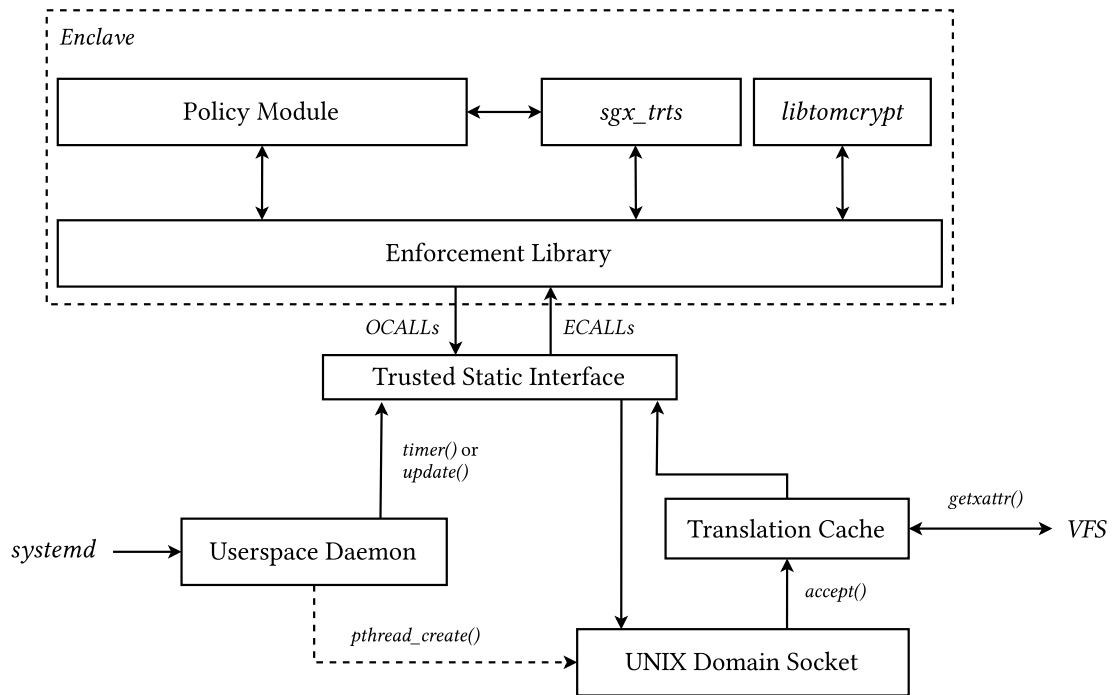


Figure 4.1: Abstract overview of an SGX enclave's protections.

Chapter 5

Evaluation

Chapter 6

Summary and Conclusions

Bibliography