

SGX Project

Team Awesome

August 15, 2016

Abstract

Web Servers are now often deployed in the cloud environment. To secure communications in the middle of the network between clients and web servers, the Internet community has come to rely on the SSL/TLS protocol which requires a strong secret (private key) to be stored with the server. However, the server administrator, in this model, is forced to trust the cloud provider to not disclose their private key. We propose a modification to legacy web servers and cryptographic libraries which secures the sensitive key material in the face a malicious provider or a compromised operating system without the need to trust the cloud provider and/or operating system.

Contents

1	Introduction	4
2	Background	6
2.1	The Principle Of Least Privilege	6
2.2	An overview of Intel®SGX	8
2.3	Threat model	9
2.4	Provisioning the enclave with the long-term private key	9
2.5	Inter-platform attestation and secret provisioning	10
2.6	SSL/TLS Overview	12
3	Design	13
4	Implementation	16
5	Project Management	17
6	Conclusion	18

1 Introduction

With the recent surge in privacy concerns, employing SSL/TLS to secure communications in the middle of the network has become common place. SSL/TLS offers guarantees of confidentiality and integrity provided that a private key's secrecy is maintained. Yet, SSL/TLS was designed assuming that its user trusts the hardware and OS of the machine on which the key is held.

While this assumption is perfectly valid in the case where a person is running an SSL/TLS enabled service on their own machines, many web applications are now hosted by third party cloud service providers such as Amazon Web Services, Heroku, Digital Ocean &c. Moreover, to offer SSL/TLS, the private key must also be stored with the web application on these service providers' machines. This implies that a server administrator using the aforementioned services is trusting the cloud-provider, including any personnel with physical/administrative access to the machines, and the underlying OS to maintain the secrecy of the sensitive key material. Such a Wide trust surface makes it difficult to maintain the privacy of critical secrets.

Consider a case where the cloud provider is not malicious; a vulnerability within their platform could lead to leaking the private key, if exploited by an adversary. Moreover, if the cloud provider is indeed malicious they could simply read your private key from the hard disk, if it is not encrypted, or mount some form of memory sniffing attack to read the key from the web server's memory since data in RAM is not encrypted. A compromised private key allows an adversary to do the following:

- Decrypt past, stored, communication between the web server and a client (assuming a cipher that does not provide perfect forward secrecy is in use)
- Decrypt any ongoing communication between the web server and a client
- Masquerade as the server and fool a client into disclosing sensitive information such as passwords

In all cases, a compromised key voids the confidentiality and integrity guarantees of SSL/TLS.

Our project aimed to break this assumption by refactoring legacy web servers to secure the long term private key in the face of: (1) An adversary who is capable of exploiting the server application, (2) a malicious cloud provider, and (3) an adversary with an exploit for the underlying operating system.

The first of these goals has been extensively addressed in previous works [1, 2] through use of the principle of least privilege - the notion that a process should only be allowed access to the smallest possible data set while still maintaining its functionality - to isolate the private key containing process from the network facing process. This approach thwarts an adversary capable of exploiting only the network facing process, as long as the interface to the private key containing process is well defined. However, it offers no protection against an adversary capable of exploiting the operating system on which the web server is running or an adversary who controls the physical machine on which the server is hosted.

Securing the private-key against such an adversary required specialized hardware called a trusted platform module (TPM), a device that can secure the private key through encryption via a Storage Root Key (SRK) [3]. The SRK's integrity is maintained by ensuring that its private component may never leave the TPM. As a result, the long term private key itself can never be decrypted outside the TPM. The TPM is also capable of executing cryptographic operations, including those that make use of the long term private key in SSL/TLS. Consequentially, by delegating all private key operations to the TPM one can rest assured that their private key cannot be compromised without compromising the TPM itself.

The security benefits of a TPM were outweighed by the cost of purchasing the additional piece of hardware, and TPMs did not gain any traction with cloud providers. In this project we utilized a new technology from Intel[®] called Software Guard Extensions (SGX). SGX is an augmentation to Intel[®]'s ISA which offer the ability to launch encrypted regions of memory, called enclaves, where only trusted regions of code can read/write. This allows for TPM like functionality, but SGX has the advantage of being deployed as part of new CPUs released by Intel[®] and as a result, when cloud providers upgrade their machines they will possess the ability to support SGX programs without purchasing additional hardware.

SGX has gained momentum as a research platform for security related work such as Haven [4] which secures a legacy application from an untrusted OS and cloud-provider *without* modifying the application's source code. Yet, there is no work, to our knowledge, that attempts to secure only the private key material through use of SGX. Narrowing the trusted region to be the component that handles the private key allows us to define a much smaller trusted computing base that only contains the CPU and the code that handles the long term private key.

The remainder of this report is divided as follows: Section 2 provides a more detailed overview of previous work and technologies underlying our project, Section 3 discusses the design of the system that we implemented to meet the above-stated goals, Section 4 highlights a few implementation considerations in realizing the system that we designed, Section 5 details the managerial aspects of this project, and, finally, Section 6 offers areas where this project could be improved, and concludes this report.

2 Background

2.1 The Principle Of Least Privilege

As previously stated the principle of least privilege (PoLP) requires that components of a system operate with the minimum resources required to complete their respective tasks. PoLP is heavily utilized in systems security research to design systems that maintain their integrity in the face of an adversary capable of exploiting some of their components. In contrast, the most popular web servers today execute as monolithic applications, including Apache and NGINX, where all of their processes have the same level of privilege and have access to the sensitive key material. Exploiting any one of these processes may therefore lead to leaking the private key.

To motivate the design of our system, we begin by discussing a system called Wedge [1] that was used to refactor Apache, by enforcing PoLP, securing the private key in the face of an adversary who is capable of exploiting the web server application. First, we list the attacks possible against a monolithic web server and then detail how the design implemented with Wedge resolved these vulnerabilities. We only consider a threat model where the adversary is capable of passively eavesdropping on secure communication channels. We did not consider the second threat model, detailed in the Wedge paper, where an attacker is capable of actively modifying packets in flight between the web server and client because the solution presented in the Wedge paper does not hold if we remove the assumption that the OS and cloud provider are trusted as the web server’s user-level processes (dynamic content generation scripts, databases &c.) need to parse web request data in the clear, and as a consequence, even if we secure the encrypt/decrypt interface using the solution in the Wedge paper, the data would be available in the non-privileged process’s memory. This is further discussed in Section 3.

Possible Attacks

There are two main attacks that may be mounted by an adversary capable of exploiting *only* the network facing component of the web server application and passively eavesdropping on packets exchanged between the server and the client:

1. The adversary could leak the private key from the network facing component, which has to run as root to bind to port 80. This could be by reading the private key from disk directly or from the process’s memory space. The attack is illustrated in Figure 1.

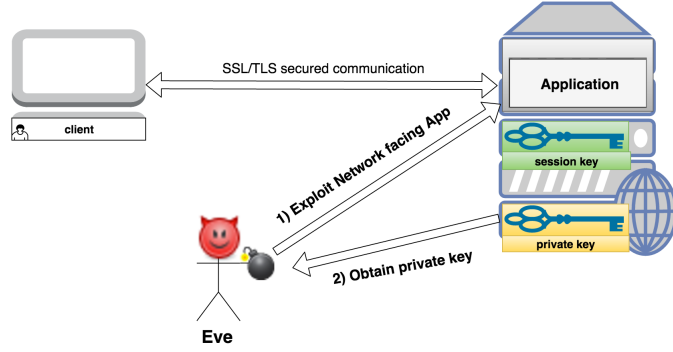


Figure 1: Exploiting the network facing component to leak the private key

2. The adversary may record traffic exchanged over the SSL/TLS channel and then exploit a naive session key generation interface to acquire the session keys used in that exchange. This exploit only works if the cipher used is one that does not offer perfect forward secrecy (PFS). The session key generation operations require use of the private key, but the adversary does not need to learn the key to complete this attack. However, the attack only affects the single session which was eavesdropped. The attack is illustrated in Figure 2.

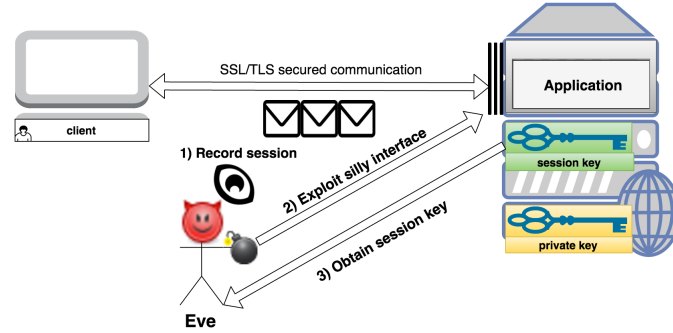


Figure 2: Exploiting the network facing component and the naive session key generation interface to generate session keys for eavesdropped session

Proposed Solution

The solution proposed in the Wedge paper is through partitioning the session key generation code into its own logical compartment¹ that executes at a high privilege level. This partitioning can be seen in Figure 3.

¹ This compartment called an sthread in Wedge's nomenclature

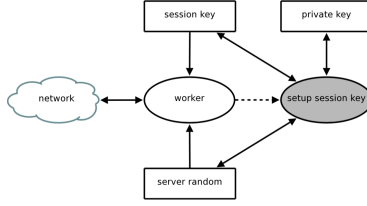


Figure 3: Partitioning Scheme to protect against leaking the private key [1]

2.2 An overview of Intel®SGX

This section will not cover all of the details of SGX but only those applicable to our project, for a complete treatment of SGX please refer to [5]. Intel®SGX is a set of x86 instructions that allow for a programming model wherein a program can be split into two components: an untrusted component that executes as normal and a trusted component that executes within a protected area of RAM, called an enclave, which can only be accessed when executing the trusted component.

The protection of an enclave is managed by the CPU; any data written to the enclave is encrypted first by a memory encryption engine (implemented in hardware) and is only decrypted when required by the CPU during the execution of the trusted component for which that enclave belongs. The key used for this encryption process is derived from a combination of a device key, unique to each SGX-enabled CPU and the “identity” of the enclave called MRENCLAVE, a cryptographic hash of the enclave’s contents at the trusted component’s initialization. SGX thus ensures that no process other than the one that initialized the enclave can access the protected area.

Interacting with the trusted component, as a result, may only occur through invoking a programmer defined interface, called a callgate, as depicted in Figure 4.

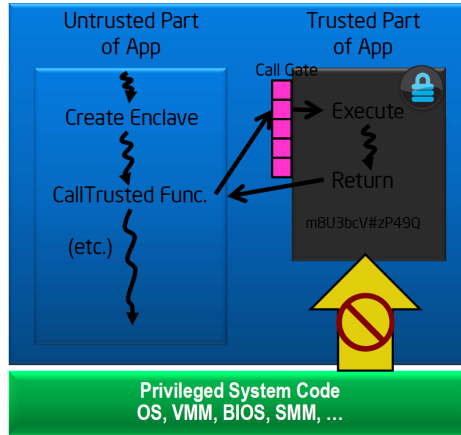


Figure 4: Interaction of the untrusted part of the application with the trusted part can only occur through a callgate

2.3 Threat model

We consider a threat model similar to the one presented in Haven [4]. Our TCB includes a correctly implemented and SGX enabled CPU and all instructions executing and data resident within an enclave. Therefore, we assume an attacker cannot access the SGX processor key provisioned by Intel itself, which is used to generate subsequent cryptographic keys that preserve the confidentiality, integrity and authenticity of an enclave.

An adversary can take full control of everything beyond the TCB. That is, we assume all software executing on the platform (outside of the enclave), the underlying operating system, the hypervisor, all firmware and the BIOS are potentially compromised. Side-channel attacks that originate from other sources such as CPU cache timing information (L1, L2, L3), power consumption or other entropy source are considered as out of scope of this work. Finally, we assume an attacker may act as man-in-the-middle to eavesdrop active sessions and, as further analyzed in the Limitations section, launch a denial-of-service attack, though without compromising any secret isolation guarantees of our design.

2.4 Provisioning the enclave with the long-term private key

Provisioning a web server with a long-term private key for the purposes of SSL/TLS is currently done storing the private-key along with the executable on the remote machine. This scheme, however, assumes a trusted cloud-provider/OS. However, under our threat model this scheme is not viable. Alternatively, we require a method by which we can verify the identity and integrity of the server application and then, upon successful verification, we can send the long term private key to the server in a secure fashion. Loosely, the requirements are as follows:

- The mechanism allows the verification of the identity and integrity of the server application and the underlying TCB. This is so that we can be sure the private-key is being sent to the same server we placed on the remote machine, and the software is being executed by trustworthy hardware.
- The mechanism allows us to setup a secure channel, ensuring that the only entities privy to the private key are the server application and the server administrator.
- The mechanism allows for the verification of the entity providing the private key. If this requirement were not in place, the mechanism could allow any arbitrary entity to authenticate with the server and provide their own private key. Observe that such an attack does not compromise the security of the long-term secret, but makes it possible to render the server useless (if the matching public certificate is not placed onto the server, verifying the server's hostname, then clients will reject connections to the server under the SSL/TLS protocol).

The first and second requirement are met by a process called inter-platform attestation, outlined by Intel in [5] and summarized in the following section.

2.5 Inter-platform attestation and secret provisioning

Inter-platform attestation is a mechanism that can be invoked by an entity, referred to as the challenger, running on one platform to verify an enclave running on another, remote, platform. This process enables the challenger to verify the following about the remote enclave:

1. The contents of the enclave's pages (code, data, stack and heap) upon creation (after the ECREATE instruction completes)
2. The identity of the entity that signed the enclave
3. The trustworthiness of the underlying hardware
4. Authenticity and integrity of any data generated by the enclave and sent as part of the attestation process. This allows us to satisfy the second requirement by generating an ephemeral key pair and binding it to the remote attestation process. This, therefore, allows the challenger to verify the integrity of the ephemeral public key and verify that it was generated by the server application.

The steps involved in the attestation process are as follows (illustrated in Figure 5):

1. The challenger invokes the remote attestation mechanism to verify the identity and integrity of the remote enclave

2. The non-trusted part of the web server receives the challenge, passes it along to the trusted portion of the web server along with the identity of the quoting enclave. The quoting enclave is a special enclave provided by Intel as part of the SGX platform to enable remote attestation by verifying the integrity of the underlying hardware.
3. The enclave invokes EREPORT which is an SGX instruction that generates a REPORT structure to be provided to a *local* enclave, the quoting enclave in this case. This structure contains a hash of the contents of the enclave's pages upon ECREATE's termination, a hash of the identity of the enclave's signer, a hash of any user-data, the ephemeral key in our case, generated by the enclave. The REPORT is signed by a MAC-key that can only be accessed by the CPU and the quoting enclave. The REPORT along with the ephemeral key is then sent to the non-trusted part of the application.
4. The REPORT is sent to the quoting enclave where its integrity is verified by calculating the MAC across its contents.
5. Assuming the REPORT is verified successfully, the quoting enclave generates a QUOTE structure that includes the REPORT structure and a signature across the quote generated using a key known as the EPID key. The EPID key is a private key unique to the CPU that is part of the platform and verifies the firmware of the processor and its SGX capabilities.
6. The QUOTE is sent along with the ephemeral key to the challenger
7. The challenger verifies the QUOTE structure by using an EPID public certificate. If this is successful then the challenger is sure that this QUOTE came from a valid SGX CPU and can trust its authenticity. The challenger can then check the contents of the REPORT contained within the QUOTE to verify the identity of the remote enclave, and the integrity of the ephemeral key received along with the QUOTE. The ephemeral key, if proven to be valid, can now be used to communicate with the remote enclave in a secure manner.

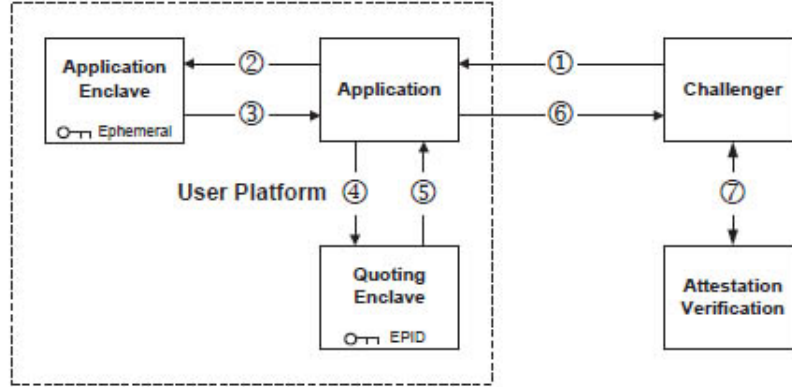


Figure 5: Remote Attestation and Secret Provisioning

However, the process outlined above takes no steps to verify the challenging entity to the trusted component. In an effort to authenticate the challenger, we utilize a combination of asymmetric cryptography and SGX’s guarantees. First, the server-administrator, before shipping the server program to the cloud-provider, stores the public key for the challenger within the text-area of the trusted component. The reason for this comes from realizing that doing so binds the public key to the identity of the trusted component, and as a result, the SGX-enabled CPU would not start the trusted-component if the public key has been tampered with. All that remains now is for challenger to sign the long term private key with their own key, enabling the trusted-component to verify it using the challenger’s public-key that was shipped along with the server-program.

2.6 SSL/TLS Overview

RSA Handshake TBC...

ECDHE Handshake TBC..

3 Design

By utilizing the aforementioned remote-attestation process, we can provision the trusted component of the remote server application with an SSL private key while making it extremely difficult* for even *privileged* processes running on the cloud provider to access the key. This is a stark contrast to current schemes wherein the private key is merely shipped as part of the server application to the cloud provider.

As previously highlighted in Section 2.6, the private key is required by a subset of the operations executed during the handshake step. These operations have to be executed within the trusted component, and are invoked by the non-trusted component via an interface that we define in this section. Note that the interface has to be carefully designed, allowing the handshake to complete correctly while not exposing the long term private key through an oracle.

Take for example an interface where the non-trusted component supplies $(\text{ServerRandom}, \text{ClientRandom}, \{\text{PremasterSecret}\}_K)$. Such an interface, while maintaining the secrecy of the private key's bits, would allow an adversary capable of exploiting the non-trusted component to generate the symmetric keys for previously eavesdropped sessions. Therefore, it is no better than leaving the private key in the non-trusted component. However, observe that it is not necessary for **ServerRandom** to be provided by the non-trusted component, it need only be provided by the *server*.

We can adjust the interface so that the non-trusted component supplies $(\text{ClientRandom}, \{\text{PremasterSecret}\}_K)$, both of which are generated by the client, and the trusted component generates a new **ServerRandom** every time the interface is invoked. The resulting interface ensures that, even if a previously eavesdropped $\{\text{PremasterSecret}\}_K$ is provided, a fresh session-key is computed on every invocation.

We consider two scenarios:

- Session keys available to the OS
- Session keys hidden inside the enclave and accessible through encrypt/decrypt oracle

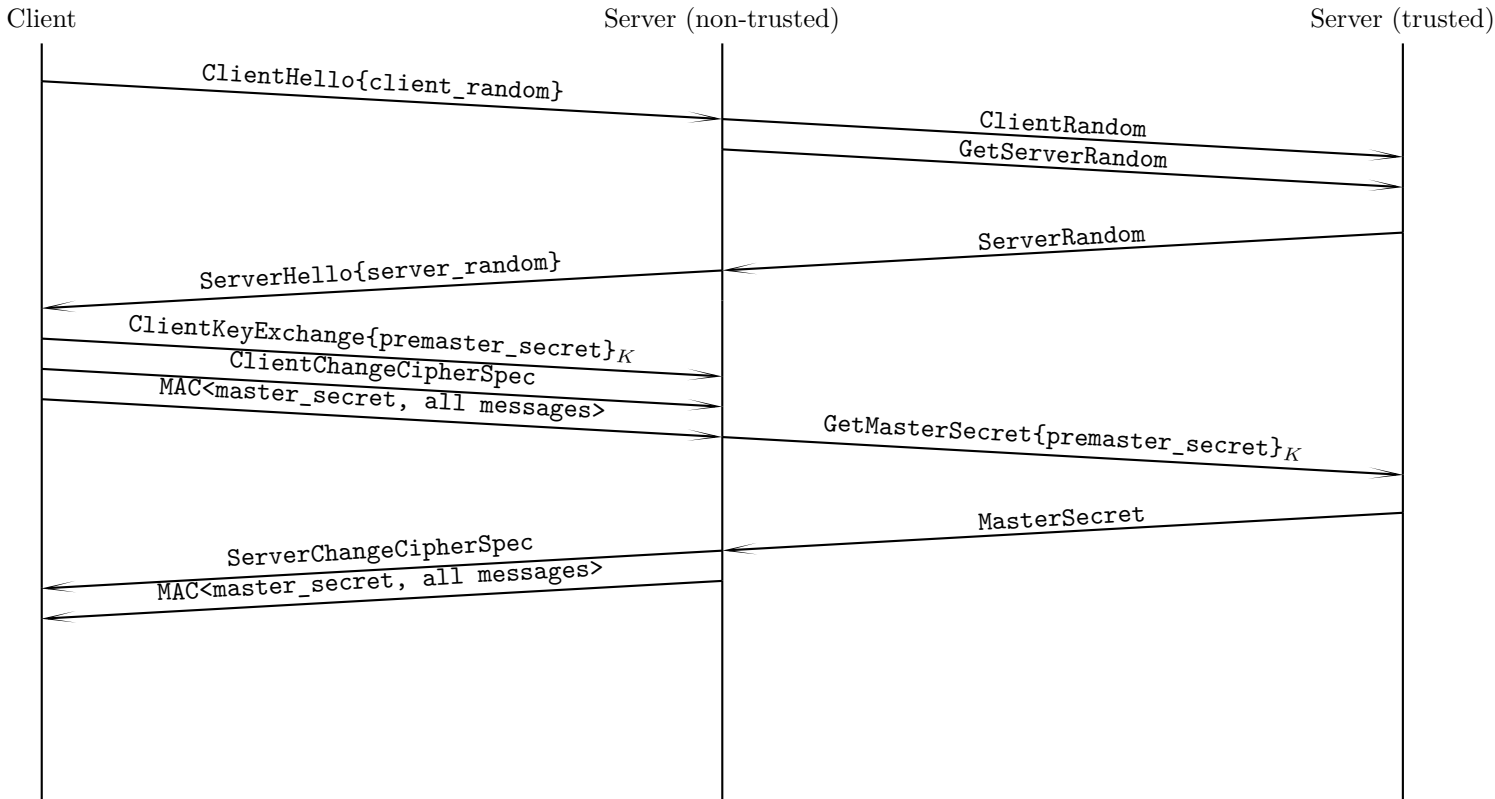


Figure 6: SSL handshake

Session keys available to the OS If our sole goal is to protect the long term, private key then we may decide it is safe to keep the active sessions keys outside of the TCB. This allows the compromised OS to use the keys as it pleases - i.e. dump all of them and send out to a man in the middle, outside of the compromised machine. However, this only permits the leakage of current (and cached) sessions keys.

One could prefer this option for performance reasons - such design puts less burden on the enclave program. In particular the role of the enclave can finish after the master secret is computed from server and client randoms and pre-master secret. This is because a pseudo-random function (PRF) lies between master secret and the input values so no information can be learned about the private key. The derivation of session keys and encryption / decryption of messages can be left unaltered in the untrusted program.

Session keys hidden within the enclave One may want to further protect the key material by requiring that the session keys themselves do not leave the

enclave. Instead an encrypt/decrypt oracle interface is presented to the untrusted component. The enclave performs the cryptographic operations using stored state accessed by the `session_id`. This prevents the compromised OS from leaking the session keys to a different machine, but requires a more complicated and expensive design. If the attacker is leaking ton of data out of the server, it may be more easily noticeable than if the session keys are exported outside of the server once.

Now the symmetric cipher context and initialization values (initialization vector,) needs to be transferred to the enclave. The symmetric keys are used to encrypt the finished handshake messages. The performance cost is increased since now at least 4 context switches between untrusted and trusted processes are required for each packet processing, i.e. 2 to decrypt an incoming packet and 2 to encrypt the outgoing result after application processing.

4 Implementation

5 Project Management

Hello world!

Hello, here is some text without a meaning. This...

6 Conclusion

Hello world!

Hello, here is some text without a meaning. This...

References

- [1] Andrea Bittau et al. “Wedge: Splitting Applications into Reduced-Privilege Compartments”. In: *NSDI*. 2008.
- [2] Maxwell Krohn. “Building secure high-performance web services with OKWS”. In: *Proceedings of the annual conference on USENIX Annual Technical Conference*. 2004. URL: <http://portal.acm.org/citation.cfm?id=1247415.1247430>.
- [3] C. Latze and U. Ultes-Nitsche. *Transport Layer Security (TLS) Extensions for the Trusted Platform Module (TPM)*. RFC Draft. University of Fribourg, 2010. URL: <https://tools.ietf.org/html/draft-latze-tls-tpm-extns-02>.
- [4] Andrew Baumann, Marcus Peinado, and Galen Hunt. “Shielding Applications from an Untrusted Cloud with Haven”. In: *11th USENIX Symposium on Operating Systems Design and Implementation (OSDI 14)*. Broomfield, CO: USENIX Association, Oct. 2014, pp. 267–283. ISBN: 9781931971164. URL: <https://www.usenix.org/conference/osdi14/technical-sessions/presentation/baumann>.
- [5] Intel Corporation. *Intel® Software Guard Extensions Evaluation SDK for Windows* OS*. 2010. URL: <https://software.intel.com/sites/default/files/managed/d5/e7/Intel-SGX-SDK-Users-Guide-for-Windows-OS.pdf>.