



UNIVERSITÉ  
DE GENÈVE

FACULTY OF SCIENCE

DEPARTMENT OF COMPUTER SCIENCE

MASTER'S THESIS

---

Confidential Computing of Machine Learning using Intel SGX.

---

*Submitted By :*

Prasad Koshy Jose

18341164

*Supervisor :*

Dr Eduardo Solana

# Abstract

Prasad Koshy Jose

June 27, 2020

Motivation : Deploying applications on the cloud requires a certain trust on third party infrastructure providers and other tenants utilizing the same platform. Unfortunately, this trust has been violated multiple times in the past. Problem - The training phase of machine learning methods requires large scale processing of sensitive data sets owned by enterprises and customers. Outsourcing the training process to a cloud-based service poses significant security risks like breaches and compromised integrity of the data and the underlying model. Methods - An effective way to overcome this security issue is to employ a machine learning system which uses a hardware-assisted Trusted Execution Environment (TEE), Intel SGX, to run the computations on sensitive data sets. The risks are potentially reduced by delegating most of the trust to the hardware. Results - Multi-user data sets are securely sealed using AES-GCM, a latest cryptographic Authenticated Encryption (AE) method and is transmitted to the machine learning code residing in a secure enclave, insusceptible to any vulnerabilities of a traditional system. Conclusion - The system is significantly secure with smaller attack surfaces and no significant performance overhead.

# Contents

<b>Acknowledgements</b>	<b>5</b>
<b>List of Figures</b>	<b>7</b>
<b>List of Abbreviations</b>	<b>8</b>
<b>1 Introduction</b>	<b>9</b>
1.1 Definitions of key terminologies . . . . .	11
1.2 Structure of the Thesis . . . . .	12
<b>2 State of the Art</b>	<b>14</b>
2.1 Introduction . . . . .	14
2.2 Background Information . . . . .	14
2.2.1 Machine Learning and Neural Networks . . . . .	14
2.2.2 Confidential and Oblivious Machine Learning . . . . .	23
2.2.3 TEE Technologies . . . . .	26
2.2.4 Attacks on TEEs . . . . .	32
2.3 Conclusion . . . . .	35
<b>3 Intel SGX</b>	<b>37</b>
3.1 Introduction . . . . .	37
3.2 SGX functionalities . . . . .	37
3.2.1 Memory design . . . . .	38
3.2.2 Data Sealing . . . . .	39
3.2.3 Software Attestation . . . . .	43
3.3 Crypto Guarantees and Primitives . . . . .	48
3.3.1 Zero-Knowledge Proofs . . . . .	49
3.3.2 Confidentiality . . . . .	50

3.3.3	Key Management . . . . .	52
3.3.4	Integrity . . . . .	53
3.3.5	Freshness . . . . .	54
3.3.6	Random Generator . . . . .	55
3.3.7	Crypto Algorithms . . . . .	56
3.4	SGX Enclave Life Cycle and APIs . . . . .	57
3.4.1	Creation . . . . .	57
3.4.2	Loading . . . . .	58
3.4.3	Initialization . . . . .	58
3.4.4	Tear Down . . . . .	59
3.5	Conclusion . . . . .	59
<b>4</b>	<b>Methods and Implemenation</b>	<b>60</b>
4.1	Introduction . . . . .	60
4.2	Use Case Solutions . . . . .	61
4.3	Process Outline . . . . .	65
4.3.1	Application requirements . . . . .	67
4.4	Software design . . . . .	69
4.4.1	Crypto Module . . . . .	72
4.4.2	Machine Learning Module . . . . .	74
4.5	Enclave Design . . . . .	74
4.5.1	Enclave Logistics . . . . .	75
4.5.2	Enclave Definition Language . . . . .	75
<b>5</b>	<b>Security Considerations</b>	<b>77</b>
5.1	Threat Model . . . . .	77
5.2	Attack surfaces . . . . .	78
5.3	Windows SSPI . . . . .	79

5.4 Side Channel Attacks . . . . .	82
<b>6 Conclusion</b>	<b>83</b>

## Acknowledgements

I would like to express my deepest gratitude to my supervisor, Dr Eduardo Solana, for his support and most importantly, for providing me with different ideas and perspectives to complete this work.

## List of Figures

1	Simple example of neural network [11]	16
2	A Simple Node	16
3	Effect of adjusting weights.	17
4	Three layer neural network	17
5	Simple Back Propagation	20
6	Convolutional Neural Networks	22
7	Recurrent Neural Networks	22
8	Use case: Hospital scenario	23
9	TPM components [2]	27
10	SEV Security Model [19]	28
11	Partitioned Deep Learning Inference Algorithm.	31
12	Threat Model of software based side-channel attacks - Cache Timing Attack	35
13	Traditional vs Trusted Remote Computing - An abstract view	38
14	Intra-platform Enclave Attestation Mechanism	44
15	Inter-platform(Remote) Enclave Attestation Mechanism	46
16	Inter-Platform communication using Remote Attestation Mechanism	47
17	Symmetric Encryption	51
18	Asymmetric Encryption	51
19	CTR mode Encryption	52
20	CTR mode Decryption	52
21	Hash Message Authentication Code(HMAC)	54
22	Freshness Guarantee using Nonces	55
23	Use case: Medical Industry	62
24	Architecture - Secure remote computation.	65
25	C1 Architecture	66

26	Components of the system. . . . .	69
27	Execution Flow Diagram . . . . .	70
28	Class Diagram. . . . .	71
29	C2 - Crypto Module . . . . .	73
30	Machine Learning Enclave Code - Class Diagram . . . . .	74
31	Attack surface with and without Intel SGX. . . . .	78



## Acronyms

**AE** Authenticated Encryption. 1, 34

**CNN** Convolutional Neural Network. 20

**CPU** Central Processing Unit. 11

**FHE** Fully Homomorphic Encryption. 14

**FLD** Fisher's Linear Discriminant. 26

**LDA** Linear Discriminant Analysis. 15

**LMC** Linear Means Classifier. 26

**LOC** Lines Of Code. 9

**NN** Neural Network(s). 12

**OS** Operating System. 11

**RNN** Recurrent Neural Network. 22

**SGX** Software Guard Extensions. 10

**TCB** Trusted Code Base. 10

**TEE** Trusted Execution Environment. 1, 10

**TPM** Trusted Platform Modules. 10

**VM** Virtual Machines. 28

# Chapter 1

## 1 Introduction

The need for data security is always on the rise. With the rise of cloud technologies, security and privacy are the most significant concerns for both enterprises and consumers. Additionally, cloud computing piles on the stress as private data is stored and processed in multi-tenant hardware providers.

Traditional systems perform access control by separating software roles into multiple user modes. Privileged system software, like hyper-visors and operating systems, control that manage machine resources and other vital components of the system, the user developed applications rely not only on the correctness of the system software, i.e. that they are free of exploitable bugs, but also on the honesty of system softwares. The correctness of the system software is vastly exploited by hackers since this is very hard to verify. The main reason for this difficulty is that a software system comprises of millions of Lines Of Code (LOC) and very complex control and data flows. Most of the operating systems, including Windows and Linux, have a lot of legacy codes, which makes them susceptible to various malicious attacks. The analysis done by [1] shows that more than 10000 computer systems vulnerabilities are discovered each year. It is not a far fetched assumption that these vulnerabilities still linger in cloud technologies and that an adversary would exploit them to gain access to confidential processes and data.

The situation gets trickier when it comes to 'honesty'. With the emergence of cloud computing, many applications and systems are deployed in third party infrastructure providers. While the pros of this approach, including cost of operations and management, are many, the cons are rather dire. The main questions that arise are 1. Where is the processing unit located? 2. How many people are sharing this hardware? 3. How many have access to the data/code? The last question is quite complex. How does one trust the infrastructure provider? Even if one trusts the infrastructure provider and all of its employees with physical access or administrative access to the

machines, one should still wonder about whether the other tenants using the same platforms are malicious..

The solution is safe data processing and execution environments inside untrusted or compromised computers. They can be used for confidential data processing in public clouds where it is necessary to have protected, encrypted code and data to protect information from malicious administrators and hackers.

As software systems have such large code bases, the attack surfaces are proportionally large. Trusted Execution Environment (TEE) inverts this logic. A TEE as an isolated execution environment that provides security features such as isolated execution, integrity of applications executing with the TEE, along with confidentiality of their assets. TEEs makes it possible to run safe, smaller pieces of code on a trusted part of the processing unit. Such smaller code bases are called Trusted Code Base (TCB). As the attack surfaces are reduced, assuming there are no back doors, the confidence in the system is evidently higher. Since TEEs are mostly part of the existing hardware, this comes with its own security flaws. Just executing a TCB in isolation(enclaves) is not enough; one might need some kind of trusted hardware which is not part of the legacy software systems.

In the last quarter of 2015, Intel announced its release of commercially available machines with its proprietary Trusted Platform Modules (TPM), called the Intel Software Guard Extensions (SGX)[3]. SGX allows out of the box attestation mechanisms in conjunction with integrity checks, memory encryption and freshness guarantees. These checks help in situations where the privileged softwares like OSs and kernels are compromised. Given the extensive researches done on its implementations, SGX is still unmatched by its competitors.

This paper explores the TEE technologies, materialized with the design and evaluation of a Intel SGX based Machine Learning system representative of the fundamental confidential computing. Focusing on the security and privacy of the model and underlying data, we wish to explore the following questions:

1. How can TPMs, in specific Intel SGX help with data security?

2. What are the challenges in implementing the software system?
3. What benefits comes from confidential computing with SGX?

## 1.1 Definitions of key terminologies

This part explains the different terminologies [4] related to TEEs and its impacts on data security. This section is inspired by the background study performed in the field and also contains an abstract view of the different technologies. First and foremost, we have to distinguish TEEs, TPMs and similar technologies like Hardware Security Modules(HSM).

We have already seen what TEEs are in the above section. Other classes of hardware for specialized cryptographic purposes exist along with TEEs, like TPMs and HSMs. TEEs are fundamentally different from TPMs and HSMs, although they can be implemented alongside the other hardwares.

A TPM is a specialized hardware chip designed to provide a "hardware root of trust" by using secret keys in such a way that physically trying to open the hardware or changing it to a different motherboard in order to retrieve its secret is next to impossible or, at the least, so that that tampering is immediately evident. TPMs are not designed to provide general computational capacity. They do provide some basic computational capabilities like generation of random keys, encryption of small amounts of data, etc. This can also be done within a TEE, but it doesn't make sense as TEEs doesn't make a good hardware root of trust. On the other hand, TEEs are a better execution environment than TPMs as they are part of the primary processors.

A HSM[4] is an external physical device that specializes only in cryptographic operations, typically encrypting a plain-text with the key it holds and returning the cipher-text so the Operating System (OS) does not handle encryption keys. Like TPMs, they are designed to detect and/or make evident physical tampering, which makes them a useful tool to hold secrets in a safe place. They generally provide higher levels of protection than TEEs, but are separate modules to the main Central Processing Unit (CPU) and motherboard, accessed via bus, network, or similar.

Neural Network(s) (NN) are a set of algorithms, modeled loosely after the human brain, that are designed to recognize patterns. NNs interpret sensory data through a kind of machine perception, labeling or clustering raw input.

Confidential machine learning is a system where the confidentiality of the training and test data is secured. In this work, confidential machine learning also protects the confidentiality and integrity of the algorithm itself.

Confidentiality<sup>1</sup> preserves authorized restrictions on information access and disclosure, including means for protecting personal privacy and proprietary information.

Integrity guards against improper information modification or destruction and ensuring information non-repudiation and authenticity.

1. Data Integrity - The property that data has not been altered in an unauthorized manner. Data integrity covers data in storage, during processing, and while in transit.
2. System Integrity - The quality that a system has when it performs its intended function in an unimpaired manner, free from unauthorized manipulation of the system, whether intentional or accidental.

## 1.2 Structure of the Thesis

This report totally consists of 6 chapters. The first chapter is an introduction to this work work, which also introduced key basic terminologies that are vital to understand the underlying technologies and principles.

The second chapter is the state of the art, which studies the latest and most relevant technological progresses in machine learning and existing confidential machine learning implementations and use cases. This chapter also explains the different types of TEE technologies in the market and gives a more vertical outlook on the top four. Finally this chapter concludes with the different attacks on TEEs and associated research works, which can be potential security considerations before

---

<sup>1</sup><https://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.800-12r1.pdf>

indulging in this work.

The third chapter is a detailed background study on Intel SGX, including its functionalities, the crypto guarantees and primitives of using this technology. The chapter concludes with Intel SGX enclaves, its associated lifecycle and APIs. This chapter solely emphasizes on Intel SGX.

The fourth chapter of this report consists of the methods and implementation of the practical work. The chapter has a brief introduction supported by use cases and the solutions that can be offered by this work. Next, the an abstract view of the process is outlined, followed by the software design of the system and its individual modules. Finally, this chapter concludes with the most important part of this work, the enclave design.

The fifth chapter is a security analyses of the system as a whole. This chapter introduces the threat model and the attack surfaces of the system. This chapter also explores the authentication system used, the Windows SSPI and credential manager, and an analyses of its individual security. Finally, the chapter includes the potential exploits that can be caused by side channel attacks to the system.

The final chapter of this report offers a conclusion and recommends possible future steps that can be taken for to have a more robust and flawless system for confidential machine learning.

# Chapter 2

## 2 State of the Art

### 2.1 Introduction

In many application domains, multiple parties would benefit from combining their respective private data sets, to train machine learning models on the aggregate data and to share the benefits of the model usage. A good use case would be that multiple hospitals can pool in patient data to train a model to aid in diagnosing a disease. Machine learning stands on the principle that having more data allows the algorithm to produce a better model and performance. But building such a strong trust-based system requires a good implementation of information security tools like encryption, TEEs, etc.

There are other ways to implement this, like Fully Homomorphic Encryption (FHE) [5] are powerful cryptographic tools that can be used for privacy preserving machine learning. But many studies in the field reports large runtime overloads, which limits their implementations in real world use cases.

First, this chapter will provide some background information on basic machine learning concepts followed by confidential and oblivious multiparty computing concepts and TEE technologies(Intel SGX). The chapter will conclude with an overview of popular attacks on TEEs.

### 2.2 Background Information

#### 2.2.1 Machine Learning and Neural Networks

In the computer science world, machine learning is defined as the ability of systems to learn from data. Systems can be trained on data to make decisions. Training is a continuous process, where the system updates its learning and improves its decision-making ability with more data[10]. A novel example of a machine learning system is implemented by Hollywood producer Vinny Bruzzese,

who uses machine learning to predict movie revenues. It is to be noted that this work only focuses on neural network based machine learning and not other classical methods like KNN or K-Means, which are avoided to stay clear of redundant work.

Systems train(learn) in two ways: Supervised and Unsupervised Learning. In supervised learning, the system is given a historical data sample of inputs and known outputs, and it “learns” the relationship between the two using machine learning techniques(of which there are several). An example would be Linear Discriminant Analysis (LDA) where the system is given a set of labelled data(like vector representations of digits along with the label). On the other hand unsupervised learning includes reorganizing and enhancing the input data to find patterns in it and to place a structure on the unlabelled data. Cluster Analysis is a trivial example of unsupervised learning. Here the system tries to find a pattern among the individual data points and tries to cluster them with respect to their similarities. If viewed from the perspective of supervised training, these clusters can be the labels of the data points.

The most wide spread application of machine learning is in data science and data science focuses on predicting(one outcome) or forecasting(multiple outcomes) the future trends.

There are several machine learning techniques, as mentioned above. Most of these techniques were first introduced almost a century ago, as they couple strongly with mathematics. In the case of neural networks, the first idea of artificial neurons was presented in 1943 by neurophysiologist Warren McCulloch and mathematician Walter Pitts. In the case of K-Nearest Neighbours(KNNs), it ages back to the early 1960s. Current computing powers makes it easier to implement these systems with ease and efficiency.

A neural network is an interconnected assembly of simple processing elements, units or nodes, whose functionality is loosely based on the animal neuron. The processing ability of the network is stored in the inter unit connection strengths, or weights, obtained by a process of adaptation to, or learning from, a set of training patterns[11].

”Network” is used to refer to any system of artificial neurons. This may range from one single node to a large group of interconnected nodes. Figure 1 shows an example of a simple neural



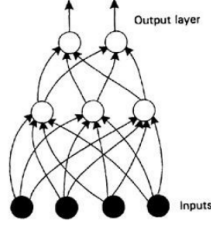


Figure 1: Simple example of neural network [11]

network. The nodes are arranged in a layered structure in which the input passes through two nodes before reaching an output. In the nodes, data is transformed using some rudimentary mathematical function. This is called a feed-forward neural network and is only one of several available. The synapses are modelled by a single number or a weight so that each input is multiplied by a weight. This product is then sent to the next layer of neurons. This is the basis of neural network based learning.

Let us consider an example of a simple node, with one input and one output. We use a simple sigmoid activation function to illustrate weights and how it affects the output.

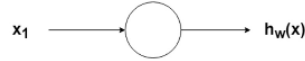


Figure 2: A Simple Node

We can observe from Figure 3 that changing the weights, changes the slope of the sigmoid activation function. This is particularly useful if there is a need to model different strengths of relationships among the input and output variables.

Neural network structures come in a myriad of different forms. The most common simple neural network structure consists of an input layer, hidden layer and an output layer(Figure 4).

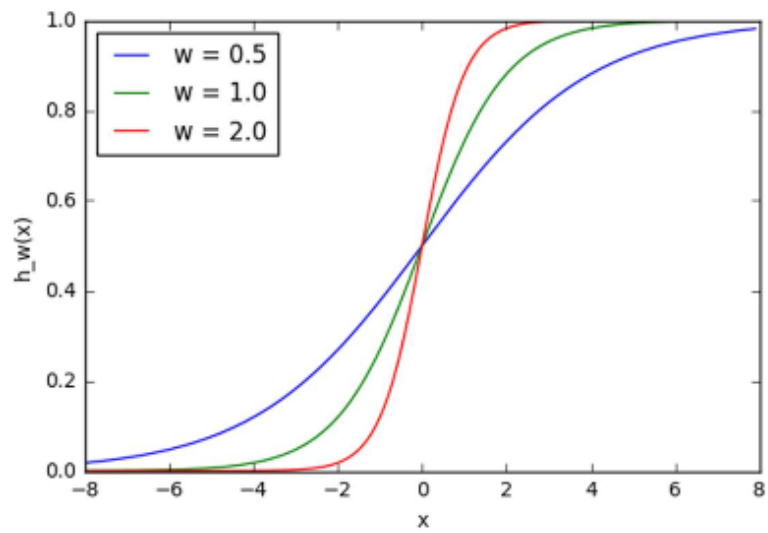


Figure 3: Effect of adjusting weights.

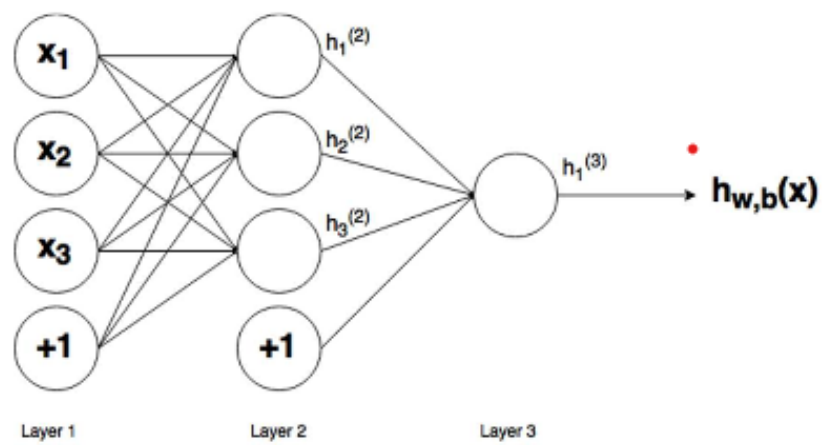


Figure 4: Three layer neural network

### Feed-forward Pass

How can one calculate the output from the input values? The above 3 layer neural network, can be used to explain the math behind it.

$$\begin{aligned}h_1^{(2)} &= f(w_{11}^{(1)}x_1 + w_{12}^{(1)}x_2 + w_{13}^{(1)}x_3) + b_1^{(1)} \\h_2^{(2)} &= f(w_{21}^{(1)}x_1 + w_{22}^{(1)}x_2 + w_{23}^{(1)}x_3) + b_2^{(1)} \\h_3^{(2)} &= f(w_{31}^{(1)}x_1 + w_{32}^{(1)}x_2 + w_{33}^{(1)}x_3) + b_3^{(1)} \\h_{W,b}(x) &= h_1^3 = f(w_{11}^{(2)}h_1^{(2)} + w_{12}^{(2)}h_2^{(2)} + w_{13}^{(2)}h_3^{(2)} + b_1^{(2)})\end{aligned}\tag{1}$$

$f(.)$  is the node activation function(sigmoid function here).  $h_1^{(2)}$  is the output of the first node in the second layer and its inputs are  $w_{11}^{(1)}x_1, w_{12}^{(1)}x_2$  and  $w_{13}^{(1)}x_3$ . They are simply added and then passed on through the activation function to calculate the output of the first node. The other 2 nodes follow similar calculation as above.

The final output of the node in the last layer is calculated from the weights( $h_1, h_2, h_3$ ) of the preceding layer and the weighted bias and not from the inputs ( $x_1, x_2, x_3$ ). Therefore, we can observe in equation form the hierarchical nature of artificial neural networks.

How can a neural network be used to classify an image of handwritten letters using this technique? The input consists of an encoding of binary data(dark and light values) of an image of one alphabet. To identify the letter class of the input, the output layer has 26 nodes, one for each letter of the alphabet. So when an input image is passed into the network one neuron is fired whenever a pattern of the corresponding class is supplied at the input.

### Back Propagation

With supervised learning, we have training data which are labelled. So while training the network, we can compare the output of the neural network to our expected training value,  $y(z)$  and feasibly look at how changing the weights of the output layer would change the cost function for the sample

(i.e. calculating the gradient). But a question arises on how to do this for all the hidden layers in between.

The method to answer this question is called back-propagation, which shares the cost function or "error" to all the weights in the network. In other words, it determines how much of the error is caused by any given weight.[12]. To calculate how much the weight has changed, say  $w_{12}^{(2)}$  has on the cost function J, we use chain function illustrated as

$$\frac{\delta J}{\delta w_{12}^{(2)}} = \frac{\delta J}{\delta h_1^{(3)}} \frac{\delta h_1^{(3)}}{\delta z_1^{(2)}} \frac{\delta z_1^{(2)}}{\delta w_{12}^{(2)}}$$

where  $h_1^{(3)} = f(z_1^{(2)})$  (sub. in eq. 1) Further solving these partial derivatives by substituting appropriate values and further applying the chain rule of derivatives, we can calculate how much of a change in weight, say  $w_{12}^{(2)}$  has on the cost function J for weights connecting to the outer layer. This can be generalized this as

$$\delta_i^{n_i} = -(y_i - h_i^{(n_i)}) \cdot f'(z_i^{(n_i)})$$

and the complete cost function derivative is :

$$\frac{\delta}{\delta W_{ij}^{(l)}} j(W, b, x, y) = h_j^{(l)} \delta^{l+1} \quad (2)$$

The output of the hidden nodes, have no such direct reference as compared to the output nodes and they are connected to the cost function only through mediating weights and potentially other layers of nodes. So the term that needs to be propagated back through is  $\delta_i^{(n_i)}$ . This is the network's best connection to the cost function. So how much does a node in the hidden layer contribute to  $\delta_i^{(n_i)}$  and how? The way a node contributes is through its weight  $w_{ij}^{(2)}$ .

As can be observed from Figure 5, the output layer  $\delta$  is communicated to the hidden node by the weight of the connection. In the case where there is only one output layer node, the generalised hidden layer  $\delta$  is defined as:

$$\delta_j^{(l)} = \delta_j^{(l+1)} w_{1j}^{(l)} f'(z_j)^{(l)}$$

In the case of multiple output nodes, the weighted sum of all communicated errors are

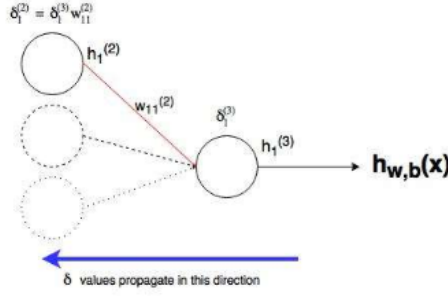


Figure 5: Simple Back Propagation

considered to calculate  $\delta_j^{(l)}$ . Hence generalised expression for the  $\delta$  values for nodes in the hidden layers is given by

$$\delta_j^{(l)} = \left( \sum_{i=1}^{S(l+1)} w_{ij}^{(l)} \delta_i^{(l+1)} \right) f'(z_j^{(l)})$$

where  $j$  is the node number in layer  $l$  and  $i$  is the node number in layer  $l+1$ . The value  $S(l+1)$  is the number of nodes in layer  $(l+1)$ . And then we calculate the complete cost function [Eq. 2].

Convolutional Neural Network (CNN) is one of the most popular deep learning methods with artificial neural networks. Neural networks are explained in the above section. This section provides an abstract view on CNNs, its working and implementation and other aspects. CNNs are widely used in the field of pattern recognition, from image processing to voice recognition. CNNs are so popular because of the lower number of parameters in ANN. Hence CNNs were widely adopted to solve complex tasks which were not possible with classic ANNs. Another important aspect of CNN is to obtain abstract features when input propagates toward the deeper layers. For example, in image classification, the edge might be detected in the first layers, and then the simpler shapes in the second layers, and then the higher-level features such as faces in the subsequent layers[13].

Convolution is defined as the process of moving the filter around the input image focusing on smaller parts called receptive fields. The first layer is the convolutional layer. For example, the top left corner of the image (say  $5 * 5 * 3$  filter), as the filter is convolving around the input image,

it is multiplying the values in the filter with the original pixel values of the image. These products are summed up into the next layer. From an abstract point of view, these filters are thought of as feature identifiers. A simple example of a feature is a curve detector. If there is a curve in the image, the sum of the products mentioned above will have a larger value than of an image without a curve. The output of this convolutional layer is an activation map. The activation map will show the areas in which there are most likely to be curves in the picture.

There are many other layers that are interspersed between the convolutional layers. A classic CNN architecture would have the structure:

*Input*  $\rightarrow$  *CL*  $\rightarrow$  *ReLU*  $\rightarrow$  *CL*  $\rightarrow$  *ReLU*  $\rightarrow$  *Pool*  $\rightarrow$  *ReLU*  $\rightarrow$  *CL*  $\rightarrow$  *ReLU*  $\rightarrow$  *Pool*  $\rightarrow$  *FC*

where CL is the Convolutional Layer and FC is the Fully Connected Layer.

For the second convolutional layer, the input is the activation maps that result from the first layer. So each layer of the input is describing the locations in the original image for where certain low-level features appear. As the network is traversed and more convolutional layers are explored, we get activation maps that represent more and more complex features. By the end of the network, we may have some filters that activate when there is handwriting in the image, filters that activate when they see pink objects, etc. The Rectified Linear Unit (ReLU) layer is used to adjust or cut-off the generated output. This layer is applied to saturate the output or limiting the generated output. The ReLU creates a sparser representation.

$$ReLU(x) = \max(0, x)$$

The pooling layer is used to downsample to reduce the complexity of further layers. Pooling layer has a similar effect to reducing the resolution. It does not affect the number of filters. A good example of pooling is max pooling where the image is partitioned to sub-regions and only the maximum value sub-region is returned to the next layer. Therefore, it should be applied only when the presence of information is important (rather than spatial information)[13].

The final layer is an important one. At this point, higher-level features are detected. The final layer takes the input volume and outputs an N-dimensional vector where N is the number of

classes the algorithm has to choose from. For the same example of digit classification, the  $N$  is 10. So each value in the  $N$ -dimensional vector is a probability of a certain class. A fully connected layer looks at which high level features most strongly correlate to a particular class and has particular weights so that the computation of the products between the weights and the previous layer will get the correct probabilities for the different classes. CNNs mainly use back-propagation to adjust

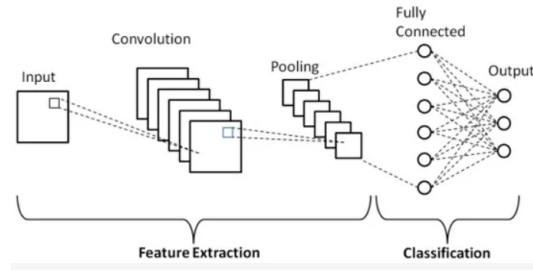


Figure 6: Convolutional Neural Networks

its filter values. Back-propagation discussed in the above section. Hence a CNN model is trained. As for the testing part, a set of images with labels are passed into the CNN and the outputs are compared with the labels.

Another famous ANN architecture is the Recurrent Neural Network (RNN). This type of ANN is widely used for learning from sequential data. The intuition behind RNNs is that the hidden state value is calculated as a combination of the previous hidden state and the latest input. Of course, each of these inputs is weighted depending on the algorithm. By always remembering the previous hidden state, RNNs are able to chain a sequence of events together. This also allows them to back-propagate errors to earlier time-steps during training. Bidirectionality can be incorporated into RNNs. Since RNNs are based on time-steps, if there was an error in any of the steps, the error

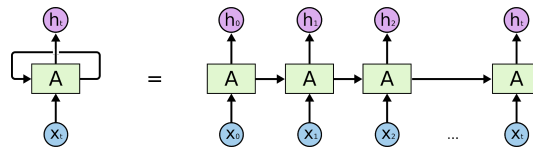


Figure 7: Recurrent Neural Networks

signal has to be propagated back several time-steps. So through each of this step, the error signal is multiplied by its layer weight and might result in a drastic effect on the magnitude of the error signal. This problem is called the vanishing gradient problem which makes it very hard to learn using a vanilla recurrent neural network. The same problem can occur when the weight is greater than one, introducing an exploding gradient, although this is slightly easier to manage thanks to a technique known as gradient clipping.

### 2.2.2 Confidential and Oblivious Machine Learning

Machine Learning consists of two stages, the training stage and the classification stage. Either one or both of these stages can be outsourced to the cloud, as illustrated by Thore et al[6]. In a traditional cloud based machine learning, there are 3 different participants involved: the Data Owner, the Cloud Provider, and the Content Providers.

The data owner is the customer of the cloud service. He owns and/or is responsible for the data that has to be processed. Content providers are devices which upload data to the cloud and which have been authorized by the data owners. Remote devices and sensors are good examples of content providers. Using the same example of a hospital system, the patient, the hospital or some large company of lab technicians are the data owners. The cloud service may be run by a third party, a partner company, a research facility or even the company itself, off-premises or in some stand-alone facility. Content providers can be the health monitoring devices like scanners which directly send data to the cloud where it is processed.

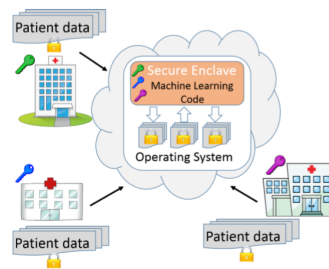


Figure 8: Use case: Hospital scenario



This kind of design is an ideal scenario for a trust based multiparty computation system. Thore et al's rationale for proposing these protocols is that, there are some scenarios where outsourcing computation to a third party cloud service makes sense from a practical and rational economic point of view. When the data is collected from multiple diverse sources, an online service can host, store and compute this data without any interaction with the data owners. These services allow the data owner to access and query their data at any time using a device of minimal computational or storage capacity. The data owner can assign different privileges to other 3rd parties interested, to access the data or receive updates concerning some other processed form of the data.

Figure 23 illustrates a model of a use case in which a medical entity, for example a pharmaceutical company, can predict the occurrence of a disease. The medical entity collects data from clinical trials, but is simultaneously interested in outside data collected from different clinics around the country and internationally. The outside clinics, meanwhile, are very cautious with their patient data and therefore are not willing to send the data outside their jurisdiction. Yet these clinics do not have the technology available locally to train the model. Through a secure confidential and oblivious platform for machine learning, the pharmaceutical company is able to easily convince the clinics to participate in the research since no raw patient data ever leaves the premises. The pharmaceutical company is also able to securely link and train their model on the combined raw data without ever seeing it.

Another viable use case is in the insurance business. In this use case an insurer who is offering business interruption and property policies needs to better comprehend the current state of a supply-chain. So if the insurer is able to focus, for instance, on a port(air, land or sea) for goods, a holistic view of the location, content and value of all relevant goods enables the insurer to better understand the risks, offer current value-in-port policies and estimate damages faster, in case of a disaster and therefore optimize the release of reserves and improve consumer interactions. In this case, the required data are within the supply-chain participants. These participants are private corporations who compete against one another and are very reluctant to release their data to the insurer as these data may negatively impact their business, if the competitors got hold of them. A

proper oblivious data sharing platform enables the insurer to give the supply-chain participants data privacy guarantees, thereby reducing the barriers to share data. Even though there is no computing needed here, this use case is of importance when it comes to data transfer to a 3rd party.

Adding to the insurance use case above, an oblivious machine learning system can be used in Multi-company fraud detection without compromising confidentiality. Churn prediction and fraud detection are standard use cases in the banking and insurance industry. As explained above in section 4.4.2, any machine learning model performance correlates strongly with the available data. In particular smaller banks often suffer from too few data points to build robust and accurate models. Currently, with this paucity of data points and lack of data science expertise and technology, the need for a trust based collaboration between such organizations is highly needed. The major catch is that banking and financial data are highly sensitive and cannot be simply shared between organizations. Hence, a confidential computing platform allows such organizations to collaboratively train machine learning models on sensitive datasets without any of the organizations having to trust each other with their data. So the more organizations participate in the system, the massively increased amount of training data can improve the fraud detection and churn prediction models. All participating organizations can use this result to optimize their existing system and hence reducing operating costs and mitigating risks.

When outsourcing computation to a 3rd party service makes sense but confidentiality of the data is an issue. The way to preserve data confidentiality is to encrypt the data before uploading it to the cloud. This may limit the usage of the data. There are many recent advances in cryptography which allow searching on the encrypted data and perform simple operations without decrypting it. A good example of one of these schemes which allows arbitrary operations on encrypted data is the FHE scheme. FHE was first constructed by Gentry[5] and it was the basis for many more subsequent schemes which have become more practical with improved performances. Since encryption involves the addition of small noise terms into cipher-texts(padding etc), operating homomorphically on cipher-texts causes the inherent noise terms to grow and correct decryption is only possible as long as these noise terms do not exceed a certain limit. The noise growth is much larger in homomorphic

multiplications than in additions. This means that FHE schemes that are in use right now can only evaluate polynomial functions of the data up to a bounded degree before the inherent noise grows too large.

The work done by Thore et al. proposes a protocol for confidential machine learning using homomorphic encryption scheme. Since many useful computational services only require evaluation of low-degree polynomials, so they can be deployed on encrypted data using only an FHE scheme or subsequent SHE or LHE schemes[[7], [8], [9]]. Experiments were conducted with Linear Means Classifier (LMC) and Fisher’s Linear Discriminant (FLD). The results showed computation time of 6 seconds with just 100 test and training data and 30 features.

The alternative to homomorphic encryption schemes is to have the machine learning service run in a trusted environment(TEEs). The following section will discuss different TEE technologies focusing on the advancements of Intel SGX.

### 2.2.3 TEE Technologies

#### TPMs

TEEs and TPMs are explained in Section 1.1. TPMs are independent and are tamper-resistance co-processors. How is the root of trust established in a TPM? The root of trust comes from asymmetric keys also known as Endorsement Keys(EKs) which are integrated and burnt physically within the module. The main functional capabilities of TPMs are also mentioned in 1.1. Some implementations also provide trusted time and monotonic counters[2]. Internal components<sup>2</sup> of a TPM is illustrated are 9.

Who can access TPMs? Access to TPMs is completely based on ownership (who sets the shared secret). When ownership is set for a TPM, a storage root key(SRK) is created and is used for sealing data etc. Only the owner is fully capable of using TPMs to its full capacity remotely, although it may be factory reset through the assertion of physical presence, either legally or illegally. Nevertheless, users other than the owner can perform operations allowed by the owner, like querying

---

<sup>2</sup>The components depends on the manufacturer of the TPM.

software measurements, storing keys and using crypto primitives.

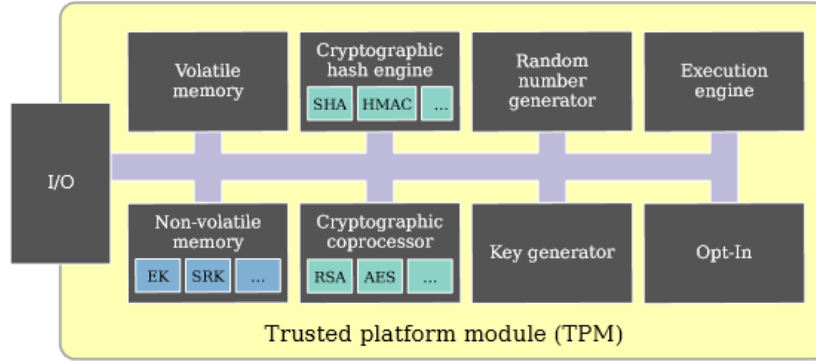


Figure 9: TPM components [2]

Being a separate component, TPMs are not a part of other isolation mechanisms such as memory page tables and privileged instructions. Security is therefore highly dependant on: 1. How TPM chips are wired to main processors, 2. What are the privileges users must hold to get access to such chips and 3. Whether adversaries might have physical access[2]. Security attacks on TPMs are analyzed in section 2.2.4.

## ARM TrustZone

ARM TrustZone[15] is a popular TEEs used in Android phones and other low energy devices. It implements a protected enclave called "Secure World" which is different from the other parts (normal world). In this architecture, each of the processors or cores, is divided into two virtual ones: secure and non-secure. The context switch is made through a monitor mode (ARMv8-A), accessed by a special instruction, or through hardware exception mechanisms (ARMv8-M). During the boot, a secure firmware initialises the platform and decides what is part of the secure and non-secure worlds by configuring the interrupt controller and setting up memory partitions and peripherals. Then, the processor switches to the normal world, typically yielding control to the OS(Normal) bootloader[15], [16].

ARM has not specified any software attestation mechanism or data sealing directly but many implementations assume the availability of a secure hardware key. This is possible as ARM chip vendors are allowed to add their own IPs to the Socket on Chip(SoCs). So it is evident that a root of trust cannot be established solely with the ARM TrustZone but in conjunction with a TPM that can only communicate with the secure world.

Another disadvantage of the ARM TrustZone is that it provides only one secure world. So if there are multiple secure applications, they have to time share or divide the protected environment, which makes it more vulnerable to attackers who can use one partition. This can compromise the entire secure world. Alternatively, allowing only one secure application to run in the secure world is possible but is very restrictive. Moreover, the secure world is under the control of a separate OS, which means that it can still be accessed by some (possibly compromised) system administrator[16]

### AMD Secure Encrypted Virtualisation (SEV)

AMD's SEV is designed to protect data in use by Virtual Machines (VM). It uses hardware encryption engine embedded in the Memory Controller (MC) which performs cryptographic operations. This method provides minimum performance impact and requires no changes to the application design. Key generation and management are performed in the AMD secure processor, a dedicated security subsystem based on an ARM Cortex-A5 processor physically isolated from the rest of the SoC. It contains a dedicated random-access memory (RAM), non-volatile storage in a serial peripheral interface (SPI) flash and cryptographic engines.

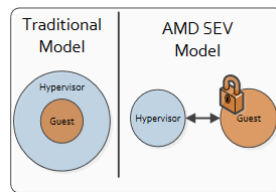


Figure 10: SEV Security Model [19]

Along with the SEV, AMD provides another TEE feature called the Secure Memory En-

encryption(SME), which provides system-wide memory encryption using a single key. In contrast to SEV, which uses multiple keys (one per VM), when using SME, ephemeral keys are randomly generated at boot time by the secure processor and are used to encrypt memory pages that are marked by system software within page tables through the C-bit[17].

SEV involves the hypervisor in many ways. The main use of the hypervisor in SEV is that it manages the keys but has to access to them. This is done by the using address space identifiers(ASIDs) and manipulating them. ASIDs are initialized and specified when the VM is launched and are used to determine the encryption key. How are ASIDs used with key management? ASIDs are used as indexes into a key table in the MC. The the key table and the MC are responsible for key generation and storage. ASIDs are also used in cache pages to tag data for faster page hits.

AMD also introduced an updated version of SEV called AMD Secure Encrypted Virtualization-Encrypted State (SEV-ES)[18]. It encrypts all CPU register contents when a VM stops running. This prevents the leakage of information in CPU registers to components like the hypervisor, and can even detect malicious modifications to a CPU register state. Before SEV-ES, the processor register state of VMs was potentially exploited by hypervisors running under AMD SEV. SEV-ES allows guests to restrict access to the register state, albeit possibly sharing it with hypervisors.

Although AMD SME/SEV/SEV-ES allow memory and registers encryption and platform attestation to a finer granularity (VMs), these technologies are vulnerable to memory integrity manipulation (e.g., rowhammer[20]) and rollback attacks, i.e., when a previous state (even if encrypted and signed) is replayed back and considered genuine.

## **Intel SGX**

This section contains a brief introduction on Intel’s SGX technology and the works done by leading researchers. Intel SGX is a TEE by Intel which employs security guarantees with respect to memory freshness and integrity. SGX limits the secure execution to user level processes, which in turn allows to purge the OS out of the TCB. It also guards against memory tampering. This is done by a physical part of the processor which helps in bookkeeping. Because of this, the memory enclaves

(protected memory) are considerably smaller in size(128 MB[21]). If there is a need to interact with the OS, there is also a considerable overhead involved.

SGX aims to guard the code execution against tampering and inspection(by infected OS) and also some physical attacks. This is done by some extension to the instruction sets in both x64 and x86 CPUs. Another feature of SGX is Data Sealing (Section 1.1). Data sealing is done by combining the data with integrity tags of memory in use by enclaves. This provides authenticated encryption and confirms the integrity of the plain-text. As a consequence, computations done in the enclave cannot be seen from the outside and any modification attempts are detected, including replay of previously authenticated values.

SGX also allows remote attestation 1.1 by which the enclaves are provided with certifications and keys. With remote attestation, a shared secret that enables a secure communication channel is established and permits the remote entity to provide the enclave with encrypted secrets.

There are many different implementations with Intel SGX. Container-based virtualization has been used to isolate applications for different platforms for interoperability. Many multi-tenant environments use Linux containers for performance isolation of applications, Docker [22] for the packaging of the containers, and Docker Swarm or Kubernetes [23] for their deployment. In multi-tenant environments, Linux containers like Docker and Kubernetes have a lower resource footprint, faster startup times, and higher I/O performance compared to virtual machines (VMs) on hypervisors. SCONE[24] is a secure Linux container based mechanism uses Intel SGX to prevent the container from being vulnerable to outside attack. The work done by Sergei et al[24]. shows how SGX can be used to plug in the security vulnerabilities of traditional container-based implementations. SCONE boasts (1) a smaller TCB and (2) a low performance overhead. SCONE offers a secure C standard library interface that transparently encrypts/decrypts I/O data; to reduce the performance impact of thread synchronization and system calls within SGX enclaves. SCONE also supports user-level threading and asynchronous system calls.

Roland et al.[26], propose a new solution, TensorScone, to extend the functionalities of SCONE to support the widely-used machine learning framework, TensorFlow[25]. Their work assists

Input:	<i>fn_enc</i> <i>bn</i> <i>lbl_enc</i> <i>img_enc</i> <i>clt</i> <i>k</i>	▶ Encrypted <i>FrontNet</i> sub-model ▶ <i>BackNet</i> sub-model ▶ Encrypted model labels ▶ Encrypted image input ▶ Client key provisioning server address ▶ Number of returned predictions
--------	--	--

```

1: ##### Within SGX Enclave #####
2: function SGX_LOAD_ENC_MDL_LBL (fn_enc, lbl_enc, clt)
3:   tls ← SGX_ATTESTATION (quote, clt)
4:   model_key, img_key ← SGX_GET_KEYS (clt, tls)
5:   fn ← SGX_VERIFY_DEC (fn_enc, model_key)
6:   lbl ← SGX_VERIFY_DEC (lbl_enc, model_key)
7:   frontnet ← SGX_LOAD_DNN (fn)
8: function SGX_INF_ENC_IMG (img_enc)
9:   img ← SGX_VERIFY_DEC (img_enc, img_key)
10:  ir ← SGX_MODEL_INF (frontnet, img)
11:  return ir
12: function SGX_CLASS_MAPPING (pv_k)
13:  result ← SGX_MAPPING (lbl, pv_k)
14:  enc_result ← SGX_ENC (img_key, result)
15:  return enc_result
16:
17: ##### Out of SGX Enclave #####
18: function INITIALIZE_ENCLAVE (fn_enc, bn, lbl_enc, clt)
19:  eid ← INIT_SGX ()
20:  SGX_LOAD_ENC_MDL_LBL (eid, fn_enc, lbl_enc, clt)
21:  backnet ← LOAD_DNN (bn)
22:  return eid, backnet
23: function INF_ENC_IMG (eid, img_enc, k, clt)
24:  ir ← SGX_INF_ENC_IMG (eid, img_enc)
25:  pv ← MODEL_INF (backnet, ir)
26:  enc_result ← SGX_CLASS_MAPPING (eid, TOP (pv, k))
27:  return enc_result

```

Figure 11: Partitioned Deep Learning Inference Algorithm.



in secure execution of existing applications that were built with TensorFlow without major changes in the commodity untrusted infrastructure. It also addresses the challenges in building such a complex system, overcoming the limitations of Intel SGX. The results show good performance with lower overheads while providing strong security with lower TCB.

Without proper confidentiality protection, input data to the cloud-based machine learning system might disclose user specific, sensitive information. The work done by Zhongshu et al.[27] focuses on the life cycles of input data instances in deep learning pipelines and identifies vulnerabilities of information leaks. They introduce a concept of ternary model partitioning based on 4 security directives established beforehand. The pseudo-code algorithm[11] showcases a basic working of confidential machine learning with Intel SGX.

#### **2.2.4 Attacks on TEEs**

TEEs are vulnerable to software flaws like memory safety violations, like stack overflows or dangling pointers. Such software vulnerabilities expose the data or can potentially change the enclave behavior. Attackers may find gadgets consisting in benign pieces of code (like the libc [28]) and reuse them to perform return-oriented programming (ROP[29]). These attacks prove that smaller TCBs have smaller attack surfaces. But using tools like SCONE[24], Haven[30], Graphene-SGX[31] propose complete run times in favor of usability as they need no modification of the legacy code. Software-based fault attacks shift the threat model from an attacker with physical access to the target device to a potentially remote attacker with only local code execution. Initially, these attacks were interesting in scenarios where the attacker is unprivileged or even sand-boxed. However, with secure execution technologies, such as Intel SGX, ARM TrustZone and AMD SEV, privileged attackers must also be considered, as they are part of the corresponding threat models.

Since most of these attacks are focused on Intel SGX and its containers, this section covers different attacks on SGX. While some of these attacks are rectified by Intel, some still remain to be fixed.

Focusing on software attacks, the work done by Jo Van Bulck et al.[35] presents a micro-

architectural attack that dismantles the security objectives of SGX by exploiting a speculative execution bug in latest Intel x86 processors. They successfully leaked plain-text enclave secrets from the CPU cache. As an optimization technique, processors may not choose to execute instructions (micro-operations) sequentially. Instead, they are executed out-of-order as soon as the required execution unit and/or any source operands become available. But towards the end of the process, these micro-ops are flushed by the processor back to the in-order and discard uncommitted micro-ops using a Reorder Buffer(ROB). Transient execution attacks abuse subtle micro-architectural side-effects to breach the memory isolation barriers. Concurrently multiple researches on transient execution attacks revealed fundamental flaws in the way the latest CPUs implement out-of-order execution. A type of Spectre[36] attacks exploit the CPU's branch prediction machinery to trick a victim protection domain into speculatively executing instructions out of its intended execution path. By "poisoning" the shared branch predictor resource, an attacker can steer the victim program's execution into transient instruction sequences that dereference memory locations the victim is authorized to access but the attacker is not.

Another research on Spectre Attacks by Guoxing et al. called SgxPectre Attacks[37] shows how to steal secrets from enclaves. Since Spectre attacks generally exploit speculative execution to subvert the confidentiality and integrity of SGX enclaves, the branch prediction of the enclave code can be influenced by processes on outside of the enclave and subsequently the enclave code can be potentially altered to execute instructions that lead to observable cache-state changes. An advisory observing the cache state changes can learn secrets of the internal registers or the enclave memory and hence defeat the confidentiality of an enclave. Furthermore, this can be used to identify the authentication tags used in AE and hence defeat the integrity of the data. The results of Guoxing et al[37] successfully applied SgxPectre Attacks to steal seal keys and attestation keys from Intel signed quoting enclaves.

Another class of attacks, called Side Channel Attacks, use alternative sources of information like power, sound, electro-magnetic or time analysis to figure out the behavior of the target and finally reveal secrets. Computer systems now use complex OSs, VMs and other processes which may

potentially leave traces that can be exploited by malicious agents[32]. One of these shared resources is the cache hierarchy, notably the most targeted component by TPM side-channel attacks [33]. There are many speculations that SGX is vulnerable to cache side-channel attacks. In fact, Intel does not consider side channels as part of the SGX threat model and thus states that SGX does not provide any specific mechanisms to protect against side-channel attacks (Intel 2016). However, they also explicitly state that SGX features still impair side-channel attacks. Intel recommends using SGX enclaves to protect password managers and cryptographic keys against side channels and advertises this as a feature of SGX. Indeed, SGX does not provide special protection against microarchitectural attacks. Its focus lies on new attack vectors arising from an untrusted operating system.

Plundervolt[34] focuses on exploiting the dynamic frequency and voltage scaling features in modern processors to corrupt the integrity of the enclave computations. Plundervolt carefully controls the processor’s supply voltage during an enclave computation, inducing predictable faults within the processor package. Consequently, even Intel SGX’s memory encryption/authentication technology cannot protect against Plundervolt. The researchers behind Plundervolt managed to recover keys or induce memory safety vulnerabilities into bug free enclave codes implemented in real world implementations. Plundervolt showcases the failure of Authenticated Encryption (AE) of data and why it is not enough to safeguard the integrity of general purpose computations.

Cache attacks exploit the timing difference between the CPU cache and the main memory. They have been the most studied microarchitectural attacks for the past 20 years. The works of Kocher [40]1996; Page[41] 2002; Bernstein[42] 2005; Percival[43] 2005 exposed vital flaws of CPUs and were able to derive cryptographic secrets. The work done by [38] was first of its kind to use cache-timing attacks to recover AES secret keys in less than 10 seconds when running inside an SGX enclave. This type of attack is called a software-side channel attack, where an advisory process thread is running on the the victim machine, trying to exploit any vulnerability. Schwarz et al[39] illustrates (Fig. 12) how one infected enclave exploits the confidentiality of the other tenants co-located on the same physical machine. This attack demonstrates a major vulnerability with SGX

by running a malware as an enclave code using the SGX features to conceal itself. The malware exploits micro-architecture implementation of CPUs, similar to Spectre[36] attacks, hence allowing the malware to infer information of other processes running on the same physical system. It is to be noted that since enclaves do not share memory with other enclaves, the OS or other processes and direct flush+reload attacks on enclaves are not possible. Same goes for any DRAM based attack from any outside processes as the hardware prevents any access to the DRAM rows in the Enclave Page Cache (EPC). However, an enclave can mount a DRAM attack on other enclaves as all the enclaves are located inside the same physical EPC.

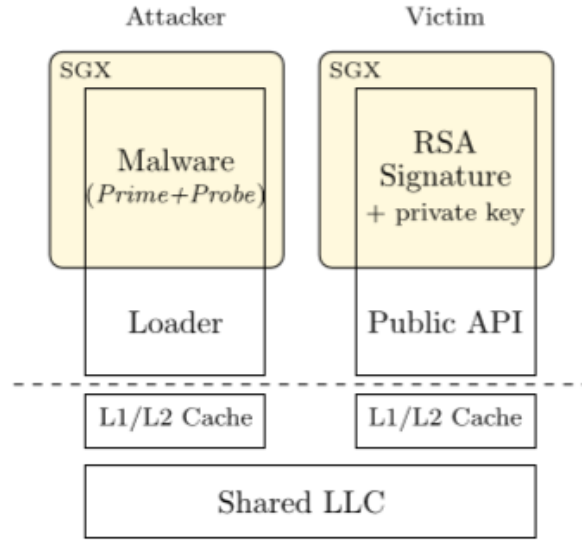


Figure 12: Threat Model of software based side-channel attacks - Cache Timing Attack

## 2.3 Conclusion

This chapter reviewed four important topics related to the thesis. The first part encompasses a brief introduction to machine learning, including the different types and techniques. Furthermore, the mathematics of two techniques of neural networks, namely feed-forward and back propagation, was explained and illustrated with examples.

The second part focused on oblivious and confidential computing and reviewed the current

trends in confidential computing and its machine learning applications, tools and runtime environments. It is vital to know how to implement the different stages of machine learning efficiently, whether to use cryptographic techniques (like HE, FHE, etc) or to use TEEs.

The third part is a study on the different TEE technologies in the market. This included TPMs and TEEs developed by leading chip makers like AMD, ARM and Intel. The different architectures and functionalities of these technologies were illustrated and studied along with their flaws. This also involved reviewing different enclave implementations, different run-times(container based) etc. Chapter 3 has the detailed working of Intel SGX as it is the technology selected for this work.

Different attacks on TEEs focusing on Intel's SGX were reviewed in the final section of this chapter. It provided a brief introduction to potential TEE vulnerabilities and attack models. This part described two main types of attacks: software based and side channel attacks. Both these attacks and leading research on them were analysed to get a view on the possible attack models and to get different perspectives of the vulnerabilities.

Hence this study shows a complete view of the current status and working of the different modules that will help in this work on confidential computing on machine learning and with Intel SGX.

## 3 Intel SGX

### 3.1 Introduction

Intel's Software Guard Extensions(SGX) is a set of extensions to the Intel processor architecture that aims to guarantee integrity and confidentiality to security sensitive processes. SGX aims at protecting processes running on a hardware where the privileged softwares(kernel, hypervisor, etc) can be potentially malicious. This vulnerability is extended to remote computations owned and maintained by an untrusted party. SGX provides some integrity and confidentiality guarantees in these situations.

The studies done on other solutions for confidential computing like FHE [2.2.2](#), shows serious performance overhead and is limited to a small family of computations. SGX aims to solve the secure remote computation problem by leveraging trusted hardware in the remote computer.

This section covers the various functionalities of Intel SGX like protected hardware and memory design, data sealing and Attestation services. The second part covers the security principles, crypto primitives fulfilled by SGX and other crypto modules built inside SGX. Finally, the SGX enclave life cycle is explained along with the versioning feature of SGX.

### 3.2 SGX functionalities

In a traditional remote computing environment, the user relied on the cloud provider, its administrators and their hardware, to perform computations on their data. Using SSL/TLS and similar tools, the user has some assurance of the computation's integrity and confidentiality. So if the cloud service provider uses a trusted hardware, it establishes a secure container and the remote computations service user can move the data and computations to the secure container. The trusted hardware, using the below explained functionalities, protects the data's confidentiality and integrity while the data undergoes computation. With trusted computing, the user trusts the manufacturer of the hardware that is in the remote computer and entrusts their data to a secure container hosted by the secure hardware. An practical use case is explained in section [23](#).

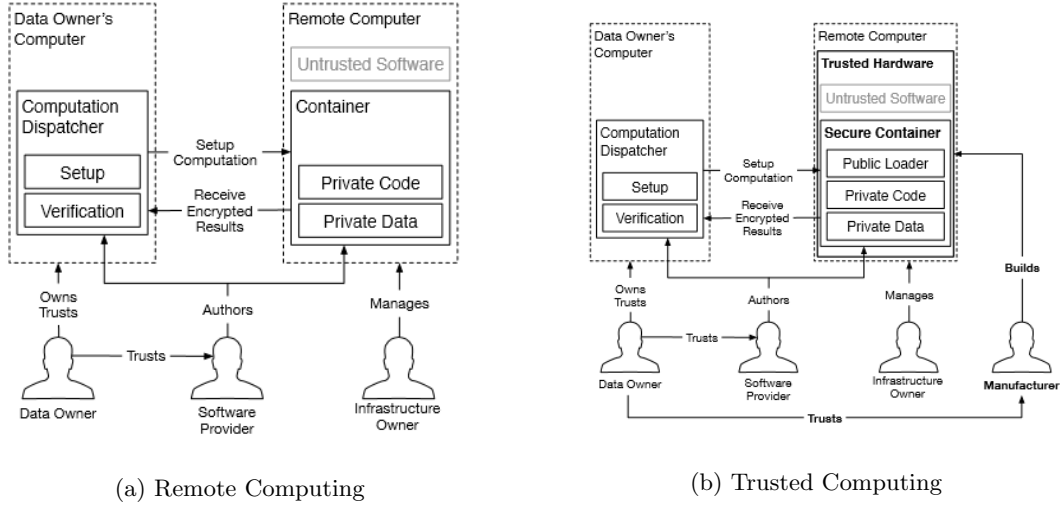


Figure 13: Traditional vs Trusted Remote Computing - An abstract view

An SGX-enabled processor that runs in the cloud protects the integrity and confidentiality of the computation inside an enclave by isolating the enclave’s code and data from the outside environment, including the operating system and hypervisor, and hardware devices attached to the system bus. This means that, the code that runs inside an enclave cannot communicate directly with the user i/o devices or the SMBs. At the same time, the SGX model remains compatible with the traditional software layering in the Intel architecture, where the OS kernel and hypervisor manage the computer’s resources.

### 3.2.1 Memory design

SGX sets aside a region in the Random Access Memory(RAM), called the Processor Reserved Memory(PRM). This region of the memory acts as the base of all enclave computations. The CPU protects this region from all non-enclave processes and memory access. As mentioned above, these restricted processes include the kernels, hypervisors and mainly the OS.

The PRM hold the Enclave Page Cache(EPC), which encompasses of 4 KB pages that stores both the code and data. Even though the EPC pages are assigned to the enclaves by system

software, this is protected, as only the EPC Metadata is disclosed to the software once the enclaves are initialized. The EPC metadata ensures each EPC page belongs to exactly one enclave. After all the enclave's pages are loaded into EPC, the system software asks the CPU to mark the enclave as initialized, at which point application software can run the code inside the enclave. After an enclave is initialized the loading method is disabled. The enclave contents are then hashed by the CPU. This is called the measurement hash, which is used by attestation services.

This aims to provide confidentiality for the code and data that runs in a remote environment as not even the hypervisor which partitions the resources to the cloud users can access a protected enclave. The PRM acts as a secure storage mechanism for the data.

### **3.2.2 Data Sealing**

Data sealing is a process of enabling data belonging to a TEE to be bound to it such a way that it can only be restored if that environment is restored. When an enclave is initialized it gives protection to the data by keeping it within the boundary of the enclave. Since the enclave memory is restricted, it is important to identify enclave secrets from others which is need to be preserved through the following events.

1. The application is done with the enclave and closes it.
2. The application is closed.
3. The platform on with the application and enclave runs is hibernated or shutdown.

What happens to the data that is in the enclave buffers after the enclave is closed? The general working involves the secrets in the enclaves to be destroyed when the enclave is closed and hence they are lost. So if there are data that needs to be persisted even after the enclave is shutdown needs to be stored outside the enclave before it is closed. Obviously this data need to be encrypted before it is stored outside the enclave. But how does this encryption work and how are the keys generated and managed? SGX has a mechanism which retrieves a key unique to the hardware



platform and also the enclave, to encrypt the data. This process of encrypting and decrypting enclave secrets is called sealing and unsealing.

Intel provides two different policies for sealing data. These are the conditions that need to be met when the data is to be unsealed.

### **Seal to the current enclave**

As we have seen in the above section [3.2.1](#), measurement hashes are generated by enclaves from hashing the enclave memory. Sealing to the current enclave uses the current version of the enclave measurement taken when the enclave was created and binds the measurement hash to the encryption key used by the sealing operation. A specific instruction EGETKEY does this binding.

Subsequently to decrypt this data, only an enclave with the same measurement hash will be able to unseal this data. The logic is that, if the sealed enclave file is tampered with, the measurement has will change as well and the data cannot be restored.

### **Seal to the enclave author**

As opposed to the above section of sealing enclave data with the enclave measurement, this policy of sealing involves using the enclave authors signature(stored in the MRSIGNER register) and binding it to the key used to seal the data. This binding is done by the hardware EGETKEY instruction. This policy not only binds the enclave authors signature, but also the ID of the enclave. Enclave IDs are generated when the enclave is instantiated[\[44\]](#). This method of sealing data with enclaves makes sure that only an enclave with the same value in the MRSIGNER register and the same ID will be able to unseal the data.

As there are two sealing policies, which one is the most secure to be used? There are two main benefits of sealing enclave data to the enclave author( $2^{nd}$  policy). First, it allows for an enclave to be updated by the enclave author without serious upgrade process or recalculating measurement hash, as it is with the first policy. Hence it is quite easier to reseal the data to a newer version. Secondly, it allows enclaves from the same enclave author to share data seamlessly.

Another subtle feature in SGX sealing to enclave author includes setting a Security Version Number(SVN) to an enclave. The CPU stores this SVN number and uses this to generate a seal key for an enclave. So an enclave has to provide this number when requesting to seal/unseal data. This makes sure that there is an option to unseal data that has been sealed by an earlier version of the enclave. This feature can be mainly used for software updates.

## Sealing Process

An abstract process of sealing data within an enclave is as follows.

1. First step is to allocate memory in the enclave to store the encrypted data and the sealed data structure. The SGX sealed data structure has the payload of the data to encrypt and the Additional Authentication Data(AAD). AAD is the additional data that will not be encrypted but for which a MAC will be calculated. AAD can be software version, headers etc.
2. The second step is to call the necessary seal data API(Section 3.3.7). The AES-GCM, Authenticated Encryption algorithm used for this work is given in 4.3
3. Finally, the seal data structure((including the key request structure) is persisted to decrypt the data.

## Unsealing Process

An abstract process of unsealing data within an enclave is as follows.

1. The first step is to allocate memory to hold the decrypted data.
2. Read the persisted key structure which was used to seal the data.
3. EGETKEY instruction is called with the key structure to obtain the seal key.
4. The unsealing algorithm API is called to decrypt the data using the above seal key <sup>3, 4</sup>

---

<sup>3</sup>It is imperative that the seal key is deleted from memory to prevent accidental leaks

<sup>4</sup>If AE is used, the hash tag generated by the decryption algorithm is matched against the tag generated during the encryption.

Two running instances of the same enclave can be distinguished at the time they attest. But there is no SGX mechanism to prevent one enclave instance to have access to the sealed data of another enclave as they both call the same EGETKEY instruction. This mechanism enables to keep data secret across power cycles[44].

A code snippet for sealing and unsealing data is as follows.

---

```
1  #include <iostream>
2  int main()
3  {
4      sgx_status_t res;
5      uint8_t* plaintext = (uint8_t*) malloc(plaintext_len);
6      // Do something in plaintext memory space.
7      // Allocate space for sealing
8      uint32_t ciph_size = sgx_calc_sealed_data_size(0, plaintext_len);
9      uint8_t* sealed = (uint8_t*) malloc(ciph_size);
10     uint32_t plain_size = plaintext_len;
11     // Seal the data
12     res = sgx_seal_data(0,
13                        NULL,
14                        plaintext_len,
15                        plaintext,
16                        ciph_size,
17                        (sgx_sealed_data_t *) sealed);
18     assert(res == SGX_SUCCESS);
19     // unseal the data
20     res = sgx_unseal_data((sgx_sealed_data_t *) sealed,
21                          NULL, NULL, plaintext, &plain_size);
22     assert (res == SGX_SUCCESS);
```

### 3.2.3 Software Attestation

Software attestation is the process of proving that a software is running on a secure and authenticated platform. In this process of demonstrating that the software has been established on a secure platform, a 3<sup>rd</sup> party establishes that the software is running inside an secure enclave on an Intel SGX enabled platform prior to provisioning that software with secrets and confidential data. The attestation process relies on the ability of a platform to demonstrate its validity by producing a credential that accurately reflects the signature of an enclave. This signature also has information about the enclave's security properties.

Intel's developer guide[44] mentions two mechanisms of attestation.

- Mechanism 1: Local Attestation between enclaves running on the same hardware(intra-platform).
- Mechanism 2: Remote Attestation between an enclave and a remote 3<sup>rd</sup> party(inter-platform).

#### Local Attestation

This mechanism of attestation is mainly used when different enclaves inside the same computer need to co-operate with one another to preform some higher-level function. Hence one enclave has to prove its identity and authenticity to the other. Initially, both enclaves create a credential, also known as a report. This report contains cryptographic proof that the enclave is in the said platform. A hardware generated key is used to get this proof. This report can be exchanges among the enclaves to prove that they exist on the same system.

This authenticated mechanism uses a symmetric key system where only the enclaves on the same platform which creates the report know the key(hardware generated) and is only embedded in the hardware platform.

A report consists of the following data.

- Measurement of the code and data in the enclave.
- Hash of the public key in the ISV certificate. The enclave receives this at initialization time.
- User data(AAD).
- Any other security related data.
- A platform verifiable signature of the above data.

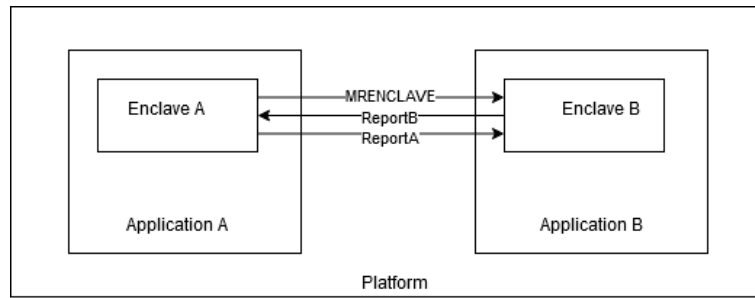


Figure 14: Intra-platform Enclave Attestation Mechanism

The figure 14 shows an example of how two enclaves can authenticate each other on the same platform.

1. Application A hosts enclave A and application B hosts enclave B. Both the applications are on the same platform. Enclave A sends its MRENCLAVE ID to enclave B.
2. Enclave B asks the hardware to produce a report destined for enclave A using the MRENCLAVE value it received from A. B transmits this report back to A.
3. Enclave A receives the report from B and requests the hardware to verify the report using its MRENCLAVE. If this is successful, it proves that enclave B is on the same platform as A. Enclave A then reciprocates this by generating a report using enclave B's MRENCLAVE value which was inside the report it received and transmits it to B.
4. Enclave B then verifies the report to confirm that enclave B exists on the same platform as A.

It is important to note that all software, including the applications outside the enclaves are untrusted. So local attestation can come in handy to verify enclaves within a single application too.

## Remote Attestation

Remote attestation proves to a user that it is communicating with a secure container hosted by a trusted platform. The proof is an attestation signature(also called a quote) produced by the platform's secret attestation key. The signature covers the container's initial state, a challenge nonce produced by the remote computer, and a message produced by the container.

An application that hosts an enclave can also ask the enclave to produce a report, which will then be used by the remote platform service to produce a quote which reflects the enclave ID and the platform state. This quote is verified using Intel(R) EPID signature verification. Hence, unlike local attestation, the CPU key is not directly exposed outside the platform.

A quote includes the following data.

- Measurement of the code and data in the enclave.
- Hash of the public key in the ISV certificate. The enclave receives this at initialization time.
- Product ID and the Security Version Number(SVN) of the enclave.
- Enclave attributes, like information concerning release or debug mode.
- User data. Similar to a [local report](#), this allows establishing a secure channel in the attestation process so that the remote server can pass secrets to the enclave that has been attested.
- A signature of the above data which is signed by the Intel EPID key.

The steps shown in figure [15](#) are as follows:

1. The challenger(cloud/service provider) issues a challenge to the application after it receives a communication request from the application which has the enclave. the challenge contains a nonce to make sure the communication is timed and live.

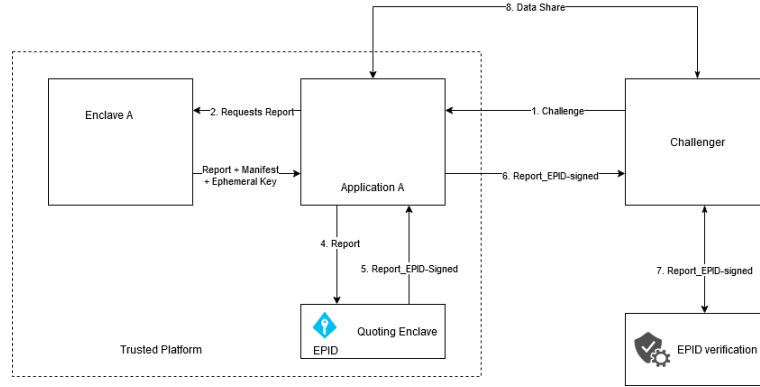


Figure 15: Inter-platform(Remote) Enclave Attestation Mechanism

2. The application requests a report from the enclave inside it. The nonce received from the challenger is sent along with the request.
3. The enclave then generates a report structure similar to the one in [Local attestation](#). This report also contains a manifest structure. A manifest can have additional information like the nonce, user data portion of the report and also a ephemeral public key in case the challenger needs to send back some secrets.
4. This report is then signed(using Intel EPID key) by the quoting enclave after it authenticates the report.
5. The application receives the quote from the quoting enclave.
6. The quote along with the manifest is sent to the challenger.
7. The challenger uses the EPID public key certificate to validate the signature of the quote. This is done with the help of a EPID verification service provided by Intel.
8. Now the challenger compared this data in the quote against the trusted configuration. If it is a match the challenger provides the necessary service to the application.

**Note:** The challenger can also add additional constraints like version number of the enclave, measurement of the data/code, the product ID, Author ID etc. This is made

for cases which sounds trivial but is a real concern. For example, an enclave developed in debug mode should never be trust with any secrets.

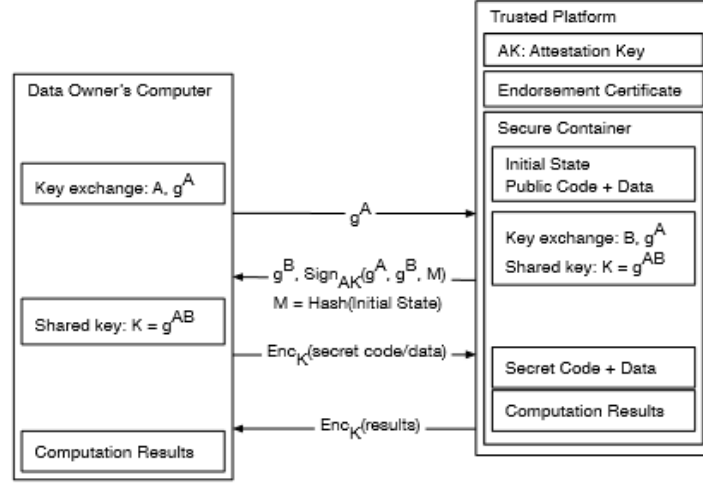


Figure 16: Inter-Platform communication using Remote Attestation Mechanism

Figure 16 puts the whole process in perspective on how the user can communicate with a remote cloud service by confirming the identity of the remote enclave, share bidirectional encrypted data.

In this use case the data owner is the challenger and their enclave is established on a trusted cloud based platform. Using the manifest explained above, where the application and the data owner exchange an ephemeral public key. After the trusted platform is attested using the EPID of the enclave, both the user and the enclave can communicate through an encrypted channel using this key and any public key cryptography.

A cloud infrastructure provider should be able to use the secure containers provided by trusted hardware to run any software on the computers. However, the provider makes money by renting their infrastructure to customers. If security savvy customers are only willing to rent containers provided by trusted hardware, and use software attestation to authenticate the containers that they use, the cloud provider will have a strong financial incentive to build the customers'



containers according to their specifications, so that the containers pass the software attestation[45].

### 3.3 Crypto Guarantees and Primitives

Systems rely on cryptographic primitives for security. These primitives have many assumptions, and building a secure system on top of them is a highly non-trivial endeavor. While the vast majority of software attacks exploit a bug in a software component, there are a few attack classes that deserve attention from architecture designers. This section covers the cryptosystems used for this work.

A message whose confidentiality is protected can be transmitted over an insecure medium without an adversary being able to obtain the information in the message. When integrity protection is used, the receiver is guaranteed to either obtain a message that was transmitted by the sender, or to notice that an attacker tampered with the message’s content.

When multiple messages get transmitted over an untrusted medium, a freshness guarantee assures the receiver that she will obtain the latest message coming from the sender, or will notice an attack. A freshness guarantee is stronger than the equivalent integrity guarantee, because the latter does not protect against replay attacks where the attacker replaces a newer message with an older message coming from the same sender.

It is also important that while signing an enclave the author chooses to remain private. Zero knowledge proofs aid in reserving this secret by a series of information sharing.

This section explains crypto guarantees and its associated primitives. Finally the respective crypto algorithms are listed and explained in detail.

<b>Guarantee</b>	<b>Primitive</b>
Confidentiality	Encryption
Integrity	MAC / Signatures
Freshness	Nonce + integrity

Table 1: Guarantees and associated Primitives

### 3.3.1 Zero-Knowledge Proofs

Direct Anonymous Attestation(DAA) is a crypto primitive which enables remote authentication of a trusted platform whilst preserving the privacy of the platform user. Intel SGX used the EPID 2.0 implementation of DAA to attest for an enclave. DAA has been adopted by the Trusted Computing Group to address this very same issue of preserving the identity of the enclave user.

The DAA protocol can be generalized based on three entities, namely, the DAA member, issuer and a verifier. The issuer has a responsibility to verify the trusted platform in the beginning and to issue(Join process) the DAA credential to the platform. The member(platform itself) uses the DAA credential from before to sign. And finally the verifier confirms the credential without violating the platform's privacy. The verifier uses zero-knowledge proofs to perform this verification of the trusted platform.

Zero-knowledge proofs provide a way for one entity to prove that it knows something without conveying anything else. Its essence is that it is not trivial to prove that one possesses knowledge of certain information not by simply revealing it, but to prove the possession of it without revealing any other information.

How does Intel implement this with its EPID? We know that Intel SGX enclaves are authenticated using the EPID cryptosystem. EPID is a group signature scheme that is intended to preserve the anonymity of the signers. Now to put in context of who the three entities are:

1. Intel's Key provisioning service - DAA Issuer Intel as the issuer is the only party that knows the issuing signing key also called the Intel Signing Key(ISK). The Group Public Key(GPK) is published within a group public key certificate signed by the ISK and is verifiable using the associated public key. This public key is different from the GPK.
2. The enclave that has to be attested - DAA member. The members join the group through an interactive protocol with the issuer. This is designed in such a way that the issuer is not aware about the member's private key. Given an enclave's report, a member creates a quote that includes an EPID signature of the report. The quote acts as an attestation that the enclave is

running on a trusted SGX platform.

3. The remote service or the user(the above use case 16 - DAA verifier. They verified that a given MRENCLAVE is running on a given trusted hardware.

The issuer publishes a group public key while securely storing the master issuing key. After an enclave authenticates itself to the provisioning service(issuer), it generates an EPID member private key. This private key is the attestation key and hence joins the group. This key is used to sign quotes.

Initially, a provisioning secret is stored in the e-fuses of each SGX enabled processor[45]. This can potentially be used by Intel to trace individual chips during the enclave authentication phase. But the join protocol can be blinded. Then the provisioning service cannot trace a signature to a attestation key and hence the chip.

### 3.3.2 Confidentiality

As mentioned above, confidentiality is attained when an adversary can intercept a message but cannot learn any information from it. This is typically achieved by encryption, a crypto primitive. The sender can encrypt the information using various techniques and the receiver at the other end can inverse this function and read the message. Most common encryption methods are block ciphers and we use block ciphers in this work. The encryption algorithms in block ciphers obfuscate the message block's content in the output, so that an adversary who does not have the decryption key cannot obtain the original message block from the encrypted output.

There are two types of encryption techniques: symmetric and asymmetric. Symmetric key encryption algorithms use the same secret key for both encryption and decryption processes. This is shown in Fig. 17.

In asymmetric key block cipher, the encryption algorithm operates on a public key and the decryption uses the corresponding private key. This is illustrated in Fig. 18.

The most popular and inherently secure block cipher based on symmetric encryption is the

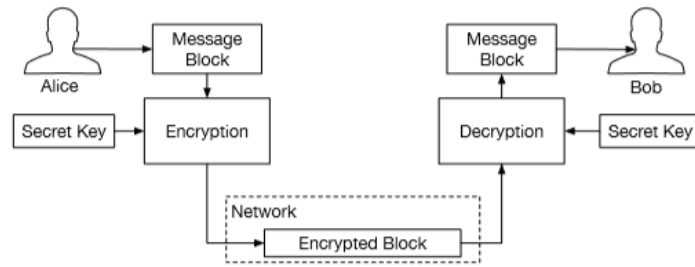


Figure 17: Symmetric Encryption

Advanced Encryption Standard(AES). AES operates on 128-bit blocks using 125 and 256 bit keys. AES transforms any 128 bit block to another 128 bit block using secure permutation. The working of AES is out of the scope of this work. Block ciphers by itself does not guarantee confidentiality completely. So it has to be combined with other operating modes. The most commonly used one is the Cipher Block Chaining(CBC). This operating mode splits the message into 128-bit blocks, encrypts them and then will be attached to the subsequent block encryption.

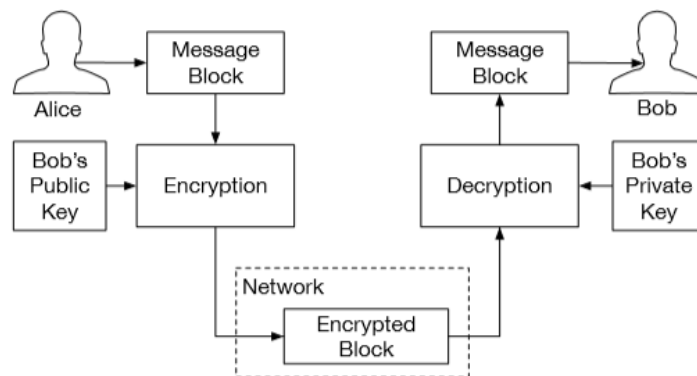


Figure 18: Asymmetric Encryption

Another example of operating mode is the counter mode(CTR). The counter mode turns a block cipher into a stream cipher. The next key-stream block is generated by encrypting successive values of a counter. A counter can be any function which produces a sequence that does not repeat for a long time. The easiest and a good example for a counter function is just an increment operator.

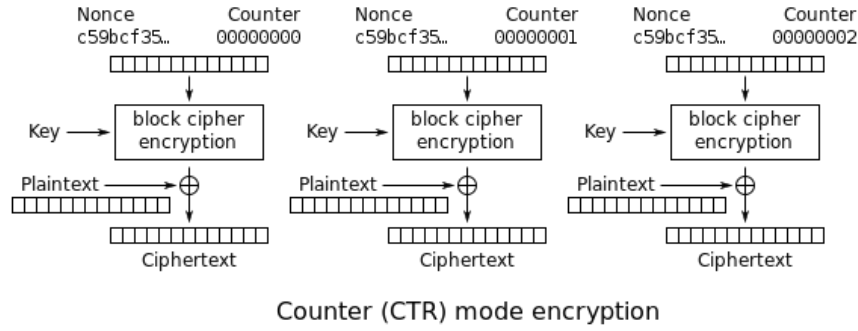


Figure 19: CTR mode Encryption

Unlike CBC which is sequential, CTR is capable to encrypt blocks in parallel. CTR modes uses an Initialization vector(IV) and a nonce for freshness. Freshness is discussed in detail in section 3.3.5. With CTR IV/nonce are combined with the counter to produce an actual unique counter block for encryption. Fig 19 and 20 <sup>5</sup> shows the encryption and decryption process of CTR block cipher mode of operation.

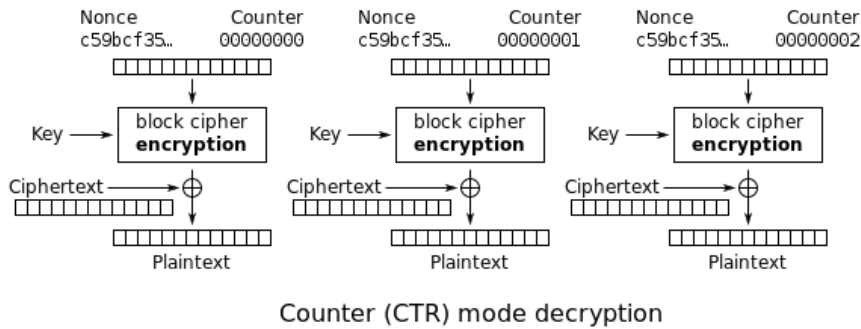


Figure 20: CTR mode Decryption

### 3.3.3 Key Management

As all cryptographic primitives rely on keys, a large part of security analysis focuses on ensuring that these keys are generated, handled and deleted according to the primitives' assumptions. The

<sup>5</sup>[https://en.wikipedia.org/wiki/Block\\_cipher\\_mode\\_of\\_operation](https://en.wikipedia.org/wiki/Block_cipher_mode_of_operation)

vital part of any key generation algorithm is the usage of random data to produce a unique key. Symmetric key cryptography requires all the parties involved in the system to establish a shared secret key. Here one party generates a secure private key using the key generation algorithm and then transmits it securely to the other parties. The channel used to distribute the key must provide confidentiality and integrity guarantees, which is a non-trivial logistical burden. The algorithms used are specified in the sections below.

### 3.3.4 Integrity

Integrity is a guarantee that a data did not succumb to any changes by a malicious third party when it was transmitted through a secure medium. Just because the data was encrypted does not guarantee that the integrity of the data had not been compromised. At the centre of integrity guarantees are the a family of functions called secure hashing functions. These hash functions operate on an *unbounded* size of data and outputs a small fixed-size output. Hash functions provides pre-image resistance which states that a malicious party cannot reverse engineer the input data from the output.

The most popular secure hash function is the Secure Hashing Algorithm(SHA)[47]. The most secure version of SHA is the 256-bit SHA-2 algorithm. This algorithm is used in this work. As mentioned above, SHA is a family of block hash functions that takes a fixed sized blocks as input. The internal state of SHA algorithms are fixed state too, meaning a fixed block size is operated upon at a time.

How can integrity be preserved with hashing functions? Similar to the encryption/decryption schemes (Section 3.3.2), integrity checks also function with a symmetric and asymmetric key setting. With the symmetric key setting, a Message Authentication Code(MAC) is calculated. A MAC of an input message is calculated with a symmetric key. The output of the MAC is fixed length and does not change for the same message and key. So at the receiving end which receives the message, the MAC and the already has the key, the MAC is recalculated and compared with the MAC that was received. The SHA algorithm is used to calculated the MAC of a message and such a cryptosystem is called a Hash Message Authentication Code(HMAC). The working of a HMAC

is illustrated in Fig 21. Here integrity is assured by computing a HMAC and transmitting it over a network along with the message. The receiver recomputes the HMAC and compares it with the version received from the sender. So if an adversary manipulated the message in transit, the receiver calculated MAC and the one received will be different.

With an symmetric key setting, integrity guarantees are provided by a crypto primitive known as signatures. Since it is an asymmetric key system, there is a private and public key. The sender provides a private key to a signing algorithm which generates a signature. The message and the signature are sent. At the receivers end, a signature verification algorithm verifies the signature using the senders public key. Signing algorithms can only operate on small messages and are computationally expensive. Therefore, in practice, the message to be transmitted is initially passed through a cryptographically strong hash function, and the hash is provided as the input to the signing algorithm.

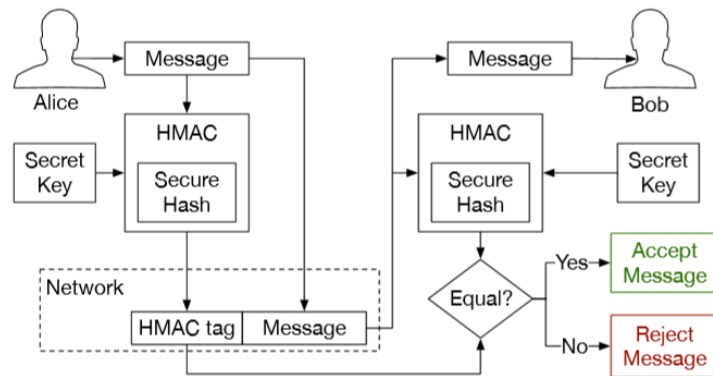


Figure 21: Hash Message Authentication Code(HMAC)

### 3.3.5 Freshness

Freshness guarantees are typically built on top of a integrity guaranteed system. This is done by adding a unique piece of information to each message. The most popular solution for freshness guarantees relies on nonce. Nonce are single-use random numbers. Nonces solve the overhead of

storing this random data on the sender side but only needs to store them on the receiver end.

A nonce is usually combined with the message timestamp and expiration scheme. This expiration immensely reduces the receivers storage requirement as expired messages can be discarded. This scheme depends on synchronization among all the entities participating in the communication. Alternatively nonces are used in challenge-response protocols for the same reason as above. The challenger generates a nonce and attaches it in the challenge request. The response to the challenge request includes an acknowledgement of the embedded nonce, so the challenger can distinguish between a legitimate new response and a replay attack.

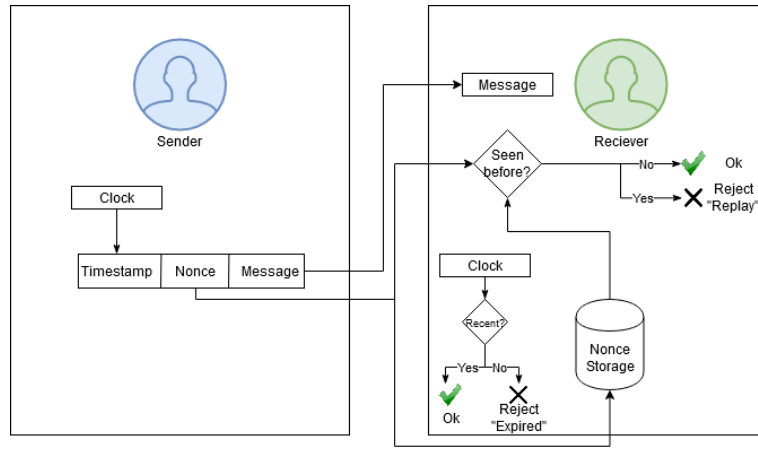


Figure 22: Freshness Guarantee using Nonces

### 3.3.6 Random Generator

SGX uses cryptographic primitives everywhere, be it to sign an enclave, to seal encrypted data off the enclave, or within the intricate software attestation mechanism. SGX also allows enclaves to use cryptography by providing a crypto libraries to enclave authors. The security of these primitives rely on a secure, random number that is used in these schemes. A good source of randomness(entropy) is necessary to build a high-quality random number generator. Commonly, in C++ uses rand() function to generate random numbers. The rand() implementation is deterministic and therefore is capable of generating predictable cryptographic secrets. For the same reasons, SGX does not



support `rand()`[44].

SGX's implementation of a random generator method is the :

---

```
1  sgx_read_rand( unsigned char* rand, size_t length_in_bytes );
```

---

This function called the hardware-based pseudo-random generation(PRNG) through the RDRAND instruction which is available in Intel CPUs. The RDRAND instruction generates a random number by provisioning hardware implementation of the underlying Digital Random Number Generation(DRNG). It is possible to call RDRAND directly but it might fail sometimes. So the SGX random generator function `sgx_read_rand()` calls RDRAND but has a retry limit of 10 times.

There is one caveat with SGX's PRNG function. As given in the SGX developer manual[44], RDRAND and RDSEED are non-SGX instructions that need to use the CPU. So an hypervisor can force a virtual machine exit. So it can prevent an enclave program from using the random function if it is malicious.

### 3.3.7 Crypto Algorithms

SGX uses the following cryptographic algorithms.

1. RSA-3072 PKCS#1 v1.5 with SHA-256, to compute enclave signatures.
2. RSA-2048 PKCS#1 v2.1(OAEP) with SHA-256, for attestation process.
3. ECDSA signatures over p256 NIST curve, with SHA 256 for enclave policy checks.
4. EC- Diffie-Hellman over p356 NIST curve for key exchange.
5. AES-GCM for data sealing
6. AES-CMAC for key derivation.
7. AES-CTR for memory encryption.

SGX only has support for AES-128 bit encryption. But RSA schemes provide a lower security of the order 112-bit for RSA-3072, and of the order 96-bit security for RSA 2048. But these

algorithms are standardized, trusted and time-tested. So if there is any security flaw, it will be in their implementations.

A good example is the cache-timing attack (Section 2.2.4) which makes the software implementations of AES very vulnerable. So to counter this, some AES implementations in critical SGX components combine hardware instructions (Advanced Encryption Standard New Instructions, AES-NI) to perform a series of rounds and a table-based implementation optimized to minimize data-leaks by this attack. In the Linux implementation of AES, a slower version of AES is used without AES-NI but other kinds of cache-timing attack mitigation.

### 3.4 SGX Enclave Life Cycle and APIs

The life cycle of an Intel SGX enclave is intertwined with resource management, specifically the allocation of the EPC pages. This section hence explains the major steps in an enclave's life cycle. Some of the major steps of an enclave life cycle are: Creation, Loading, Initialization and Tear down. It is important to study these steps to effectively work with SGX enclave development.

#### 3.4.1 Creation

There are specific instructions that handle and manage an enclave life cycle. The ECREATE instruction creates a new enclave by turning a free EPC page into a SGX Enclave Control Structure(SECS). ECREATE also validates the information used to initialize the SECS and results in a page fault or a general fault, if the information is not valid. A good example of this is to specify the enclave size as a power of 2. So if this condition is not satisfied while setting the SIZE field, then ECREATE results in general fault. Finally ECREATE initializes the enclaves INIT attribute to false. This is so that, the enclave's code cannot be executed until the INIT attribute is set to true. This happens in the initialization step.

### 3.4.2 Loading

This step in the life cycle is when the code is loaded into the enclave. When the enclave's SECS is in the uninitialized state (set in the creation step), the system software uses the EADD instruction to load the initial code and data into the secure enclave.

EADD reads input data from a Page Information structure(PAGEINFO). The PAGEINFO[45] structure can only be used to communicate information in the SGX implementation. EADD validates its inputs before modifying the newly allocated EPC pages. So if attempting to EADD an EPC page to an enclave which is in initialized state will again result in general fault. Similarly, if attempting to EADD an EPC that is already allocated results in a page fault.

This step is also responsible for the measurement calculation(MRENCLAVE) that is used in the software attestation in Section 3.2.3.

### 3.4.3 Initialization

After loading an enclave in the loading step, where the data pages are loaded into the enclave, the system software uses a Launch enclave(LE) to obtain a EINIT token structure. This token marks the enclave's SECS as initialized when passed to the EINIT instruction. The LE is a privileged enclave from Intel and is a prerequisite for the use of enclaves. Since LE is a SGX enclave it has to be created, loaded and initialized as a normal enclave. Also the LE is cryptographically signed with a key that is hard-coded into the hardware. This makes EINIT to initialize the LE without checking for a valid EINIT Token Structure.

After the initialization step is successfully, it sets the INIT attribute to true. This enables the software application to execute the enclave's code. Furthermore, once INIT is set to true, EADD cannot be invoked on that enclave anymore, so the system software must load all the pages that make up the enclave's initial state before executing the EINIT instruction. This makes sure of the enclaves integrity guarantee.

#### 3.4.4 Tear Down

Finally, once the enclaves purpose is completed, the EREMOVE instruction deallocates the EPC pages used by the enclave. EREMOVE marks an EPC page as available by setting the VALID field to 0. Also before freeing up the page, EREMOVE confirms that there is no logical processor executing the enclave code that has ownership of the page that is about to be destroyed.

An enclave is destroyed completely when the EPC page that hold the SECS(which is created in the first step) of the enclave is freed. it also will not deallocate a SECS if ti is referenced by any other encalve. This is to make sure that an enclave's SECS page can only be deallocated after all the enclave's pages have been removed.

As mentioned above it is important to understand an enclave's lifecycle to avoid any memory leaks or unreferenced or deallocated memory is removed and hence is free of any vulnerabilities.

### 3.5 Conclusion

This chapter extensively covered the basic working of Intel SGX, namely memory design, data sealing and software attestation. This knowledge comes in handy while incorporating SGX into a legacy application or developing a stand-alone SGX application. The two different software attestations and their respective steps for implementations were illustrated and explained with an example.

The second part of this chapter highlights and explains the crypto guarantees and primitives that make SGX a viable security solution. The working of each of these primitives are explained to get a different perspective of how the above mentioned functionalities can be helpful with real-life threat vectors and scenarios. This also explains the assumptions of these primitives, which makes any security system vulnerable if not implemented properly.

This chapter concludes with the lifecycle of an SGX encalve, from its initialization to tear-down. This includes all the associated sub-process and structures which aides in the development of an SGX application. Life cycle of an enclave is an important knowledge to possess to properly create and discard enclaves without any data leaks or security vulnerabilities.

## 4 Methods and Implementation

### 4.1 Introduction

Software applications often need to work with private information, in this case the data sets which can be financial information, encryption keys, health records etc. This sensitive data is intended only to be accessed by the designated recipient. In the context of this work, we refer to this as application's secrets. The OS's job is to enforce security policy on the computer system so that the secrets are not unintentionally exposed to other users or applications. Applications often employ safeguards for their data, such as data encryption, to ensure that data sent to storage or over a network connection cannot be accessed by unauthorized third parties even if the OS and hardware are compromised.

Despite these protections there is still a significant vulnerability present. While there are many safeguards in place to protect one application from another or the OS from an unauthorized or unprivileged user, an application has no protection from higher privileged processes like the kernel or the OS itself. So if a malware that gained an privileged access to the OS or the kernel, it has unrestricted access to all system resources and applications running in the system. This same problem reflects in a cloud based platform too. So a sophisticated malware which infiltrated the cloud system can target an application's protection schemes to extract encryption keys and even the secret data itself directly from memory.

Hence, we propose a TEE based application to securely upload multiparty data-sets. The application solely works with Intel SGX and will require a SGX enabled machine to run on and with the right configurations. The application is designed (Section 4.4) to reduce the attack surface and to keep it to a minimum. Furthermore, the different classes of the software are illustrated and explained in this chapter. Section 4.5 of this chapter focuses on the enclave components and the flow of application secrets among them. This chapter concludes with all the security considerations and pitfalls of the implementation and of the underlying technology.

## 4.2 Use Case Solutions

There are three use case scenarios (section 2.2.2) studied for this work. These use cases are spread among the top most data sensitive industries and which would require complete process transparency but data confidentiality. The cut throat competition among them makes it difficult to come up with a robust solution for collaborative computing or data transfer among them or with a 3rd party. This work proposes a system to provide a general viable solution to each of these use cases.

### Medical Industry

In this use case a pharmaceutical company has a machine learning model that can predict the occurrence of a disease to a certain extent. The company can then use this result to change production of medicines, Personal protective equipment (PPE), vaccines etc. The data comes from clinical trials, most commonly sponsored by the same pharmaceutical company. But they are also interested in having more data coming from clinics from a different part of the country or from a different part of the world. Now the problem here is that these clinics(independent or sponsored) are very cautious with the patient data and therefore are very reluctant to share it outside their jurisdiction. But at the same time the clinics do not possess the required hardware nor the knowledge to train the model locally at their sites.

**Solution:** A confidential machine learning platform built on top Intel SGX can potentially solve this problem. The application enables the clinics to encrypt their data locally and then send it over to the secure machine learning platform. The model gets applied to the data and computations are performed to get accurate predictions. Finally, the results are also encrypted and sent to all the participating parties. This enclave based system guaranteed that at all times the data and the model remains confidential and tamper resistance with integrity checks in place at all critical points.

Hence, the pharmaceutical company is able to

1. provide guarantees to the clinics that no raw patient data ever leaves their jurisdictions and hence can convince them to participate in the research.

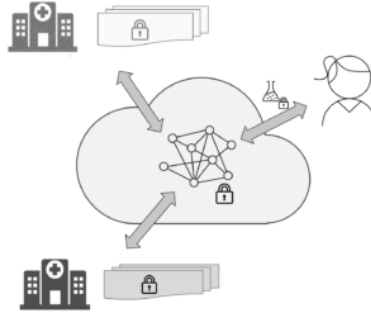


Figure 23: Use case: Medical Industry

2. establish a secure link throughout the community and train their model on all the combined data without actually having access to it.

### Insurance Industry

This use case mostly involves data sharing. Trust based data sharing is a vital component when it comes to oblivious computation systems. This use case shows that one component of this proposed work can be used to solve an important use case. The insurance industry is well known for its trust issues. A business interruption and property policies insurer has a need to better understand the current state of an entire supply chain. The insurer has to better understand the risks, current value-in-port policies or faster estimate damages in the case of disaster against which that particular business is insured against. So a holistic view on the port of the goods or business like the location, content, value of all the relevant goods and transactions will aid the insurer to make proper decisions to release reserves and improve customer interactions. The current problem is that the required data is with the supply-chain participants. These participants are reluctant to share data as they are very sensitive and if a competitor got hold of them, it would negatively impact their business. So the communication among the supply chain participants and the insurers are not effecting resulting in improper customer interactions.

**Solution:** An enclave based system solves trust issues by the process of remote attestation. After the entities undergo remote attestation, the supply chain participants can share the

cryptographic keys to decrypt their separately uploaded data-sets. The application can then be queried for the appropriate results. The database is designed for the enclave and has the capabilities to store the data sets. Queries can include the total value of goods in a port over time or distributions of time-in-port over all goods.

Hence the data remains confidential in the enclave and the insurer can provide data integrity guarantees by MRENCLAVE calculations etc.

## **Banking Industry**

Banks depend a lot on churn prediction and fraud detection to identify illegitimate transactions and customers. Churn prediction uses big data to detect customers who are likely to cancel their subscriptions, accounts and services. It is used by most of the large subscription businesses and banks to identify customers most at-risk of churning. If a good churn prediction model is in play, banks can have huge business savings and customer base, even if implemented in small scale. Fraud detection<sup>6</sup> is a set of activities undertaken to prevent money or property from being obtained through false pretenses. In banking, fraud may include forging checks or using stolen credit cards. Banking and insurance industries use standard machine learning for churn prediction and fraud detection. As with all machine learning, the more data that is used for training, the higher the prediction score is. This means that the model performance highly correlates to the available data. Banks in general, not just small banks, often suffer from very few datapoints to build good and high performance models. In the case of smaller banks, combine this issue with lack of data science infrastructure and expertise, creates a desire and need to collaborate among similar organizations.

**Solution:** Similar to the above use cases the solution relies on the trust that is developed with an enclave based application. For example, the machine learning algorithm for fraud detection runs as an enclave based service in the cloud, which is ready to receive as many data points as possible. Using SGX capabilities like software attestation(remote), the data is shared with the service after a required trust has been established. The service then runs the training and outputs

---

<sup>6</sup><https://searchsecurity.techtarget.com/definition/fraud-detection>



the resulting model to all the participating organizations in a secure way.

This allows all the participating organizations to collaboratively train machine learning models on confidential and sensitive data without any of these organizations having to trust each other. The organizations initially establish trust only with the platform and this is done in a very transparent manner. The massively increased amount of training data now results in very accurate prediction scores and high performance in detecting fraud. This in turn reduces the cost for the organizations.

### 4.3 Process Outline

The design of the system requires the application code to be divided into two main parts of code bases: The trusted code base and the Untrusted code base.

1. **Trusted code base** is the enclave itself. The code base that resides in the trusted code is the code that has access to the application secrets. An application can have more than one trusted code base or in other words, an application can have more than one enclave.
2. **Untrusted code base** is the rest of the application, its UI modules etc. In a TEE design, the OS, the kernel and the VMM are considered to be part of untrusted code base.

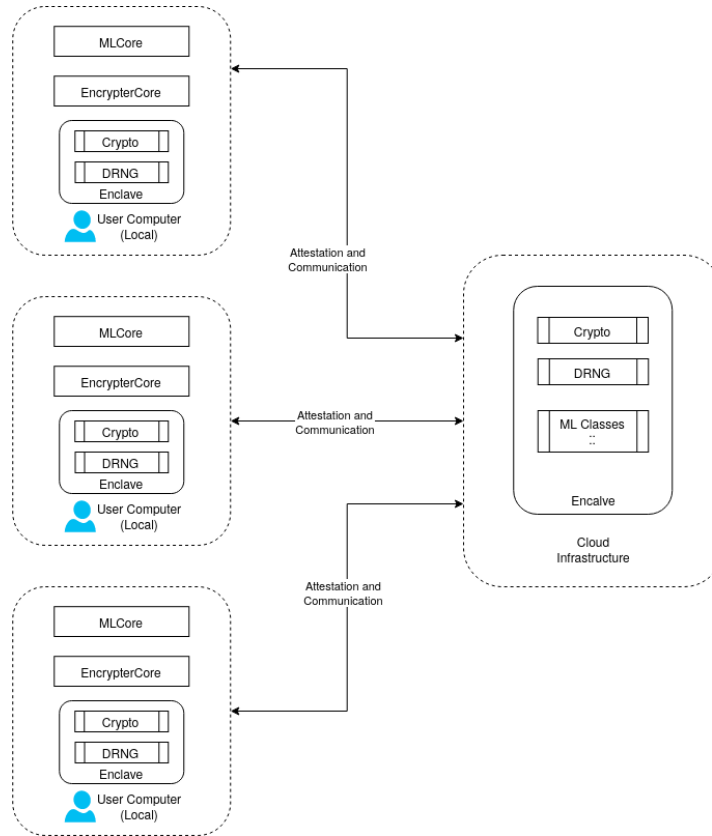


Figure 24: Architecture - Secure remote computation.

Intel SGX makes it possible to run a code on a remote system, typically a cloud provider

such that the remote system cant alter the program that is executes(maintains integrity) and also the program can process encrypted data that is decrypted only inside the CPUs(cloud) and therefore inaccessible to the OS, VMM or the kernel. The system is designed to move the trust from the cloud provider to the application code and Intel SGX.

The illustration shown in figure 24 shows a top level view of implementation. Once a secure connection and trust is established through attestation the user initiates sensitive data communication to the cloud. Since the machine learning code resides in the cloud, the data flow inside the enclave is in plaintext.

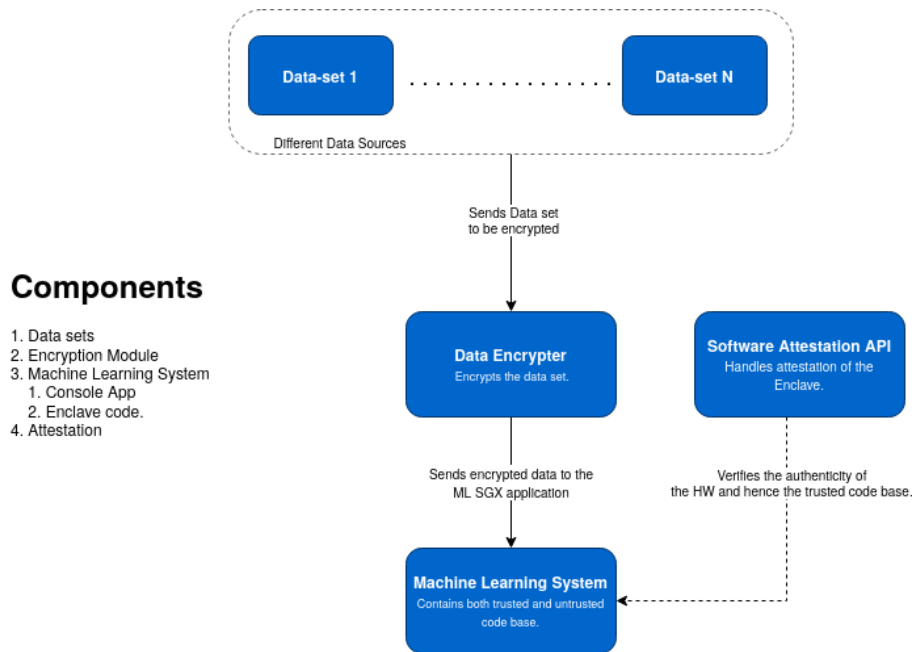


Figure 25: C1 Architecture

To preface the core software components, classes and their design, an abstract view of the architecture is given in Figure 25. The systems contains 2 main modules: the Encrypter Module and the Machine learning module. The Encrypter module is responsible to seal the data using authenticated encryption. The Encrypter module resides completely in the users system and is bound to the enclave that is used to seal the secrets.

The Machine learning module resides in a SGX enabled cloud system. This module is responsible to decrypt the application secrets and to perform the necessary computations. The machine learning module is not exposed to even the OS, the VMM or any third party. Before the data is shared, the user confirms the identity of the ML module using SGX software attestation (Section 3.2.3). Depending on the use case and implementation, this attestation can be local or remote. For this work, since the cloud environment is simulated in the local system, a local attestation is done to authenticate the identity of the enclaves before secure data communication.

Steps to use the application are as follows:

1. User registration is done and keys are generated using Windows Security Support Provider(SSP) API.
2. The user authenticates himself to the system using Windows Security Support Provider Interface(SSPI) and establishes a session. This session is timed and recorded by the local system.
3. The user encrypts the data set(s) using the encrypter module and the encrypted file(s) are saved locally.
4. The user uploads encrypted data sets to the machine learning module through the application UI.
5. The user trains on all the data in the repository and is given the learning and test accuracy.

#### **4.3.1 Application requirements**

Some basic application requirements can help narrow down the scope of the work so that the focus can solely be on the Intel SGX integration rather than minute design and development of the application.

*Requirements and  
Design Decisions*

1. Only runs on Intel SGX capable systems or Visual Studio 2017 with Intel SGX SDK.
2. Requires Microsoft Windows 10 64-bit for Intel SGX support.
3. Requires an Intel processor that supports Intel Data Protection Technology(DPT) with Secure Key
4. Must not rely on 3rd party libraries or utilities.

Table 2: Requirements and Design Decisions

---

*Requirements and  
Design Decisions*

1. C++ based Win32(WinAPI) Front end and Back end
2. Requires Microsoft Windows 10 64-bit for Intel SGX support.
3. Must not rely on 3rd party libraries or utilities.

Table 3: UI Requirements and Design Decisions

---

Table 2 shows the system requirements to run the proposed work. The 3rd requirement gives access to RDRAND instruction which ensures access to high quality entropy source. Table 3 shows the requirements of the UI module and the corresponding design decision impacts. The main reason to use Win32 for the application is to reduce usage of managed code like frameworks etc. Developing applications with managed code adds an extra layer of complexity and communication

flows which in turn has an impact on the attack surfaces.

## 4.4 Software design

The trusted component should be as small as possible, limited to the data that needs the most protection and those operations that must act directly on it. A large enclave with a complex interface doesn't just consume more protected memory: it also creates a larger attack surface.

Enclaves should also have minimal trusted-untrusted component interaction. While enclaves can leave the protected memory region and call functions in the untrusted component (through the use of a special instruction), limiting these dependencies will strengthen the enclave against attack.

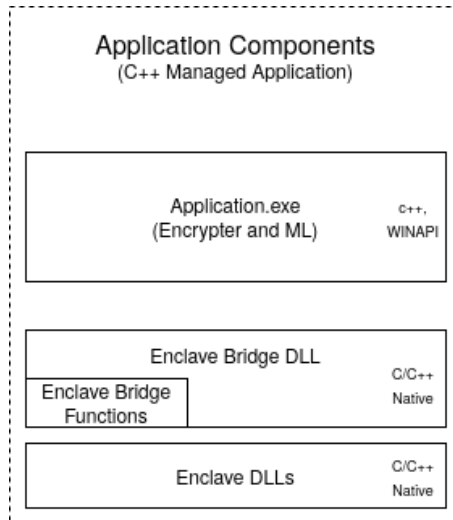


Figure 26: Components of the system.

There are two main parts of this application. One is the trusted part and the other is the untrusted part. The trusted code resides in the enclave DLLs as shown in the Figure 26. The enclave code can be C or C++ and are secured using Intel's SGX security guarantees. One main drawback of SGX is that enclaves cannot directly interact with the input/output modules of the system. This means it is not possible for a user to directly communicate with the enclave code.

How to communicate with the enclaves? This communication usually entails sharing of data, which can be application secrets, from the UI to the enclave. Even though this data can be encrypted, there are other security concerns like freshness or integrity of the data. This communication can be done effectively using bridge functions which are usually wrappers around the ECalls or OCalls. The implementation of bridge functions solely depends on the application design and the technology stack used. If the application is built using managed code, then these wrapper functions can bridge the gap between the enclave code which is written in C/C++ and the application. For this work, bridge functions are written even though they are redundant, to show how to robustly develop an Intel SGX application.

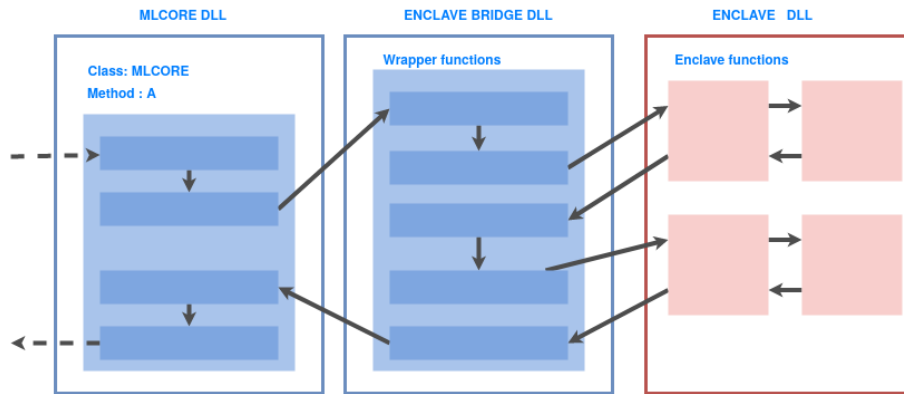


Figure 27: Execution Flow Diagram

The Figure 27 shows how an untrusted part of the application can communicate with the secure enclave code. Majority of the ECalls are simple wrappers around the class methods in vault. Enclave bridge DLL comprises of wrapper functions of the methods of the enclave classes. This is done because of the complication that ECalls, OCalls and bridge functions must be in native C code. Hence, once the enclave is launched there is a need to span the gap between C and C++ objects.

Figure 28 shows the class diagram of the system. The one path of communication between the trusted and untrusted code is from the FileHandler Class to Crypto Class. The encrypted data is decrypted once it enters the enclave and moves around as plain-text within the enclave. The Crypto methods are

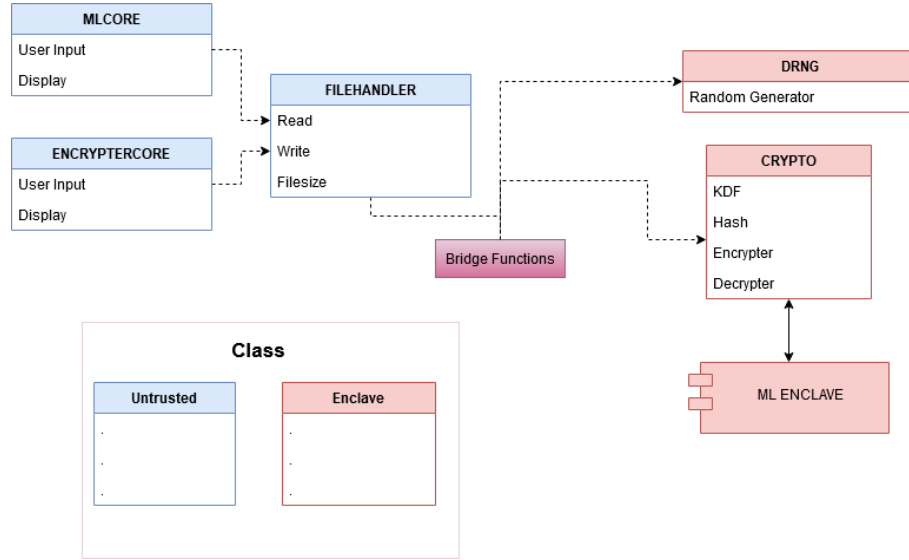


Figure 28: Class Diagram.

1. Encrypter - The encrypter method reads the data and seals it using a local enclave. The sealing function used is AES-GCM which is described in Section 3.3.7. Data can be sealed using two modes: Seal by Enclave ID or Seal by Author. For this work and simplicity, the data is sealed by using enclave author mode.
2. Decrypter - This method resides as enclave code in the cloud as part of the ML module and unseals the data passes to it. The decryption module can only unseal data after the encrypter enclave has been authenticated by the cloud enclave once per session using software attestation methods of Intel SGX.
3. KDF - The Key Definition Function generates a 128 bit key that can be used for AES using SHA.
4. Hash - Hashing method to generate secure hashes.

Another class in the enclave code is the DRNG class. The DRNG is the interface to the on-chip digital random number generator which uses the Intel Secure Key. There are SGX alternatives to the Intel Secure Key technology. The DRNG takes care of the RDRAND and RDSEED feature



detection checks and instruction calls. The DRNG class<sup>7</sup> generates 512 pairs of 128-bit values and mix the intermediate values together using the CBC-MAC mode of AES to produce a single, 128-bit seed. The process is repeated to generate as many seeds as necessary.

The MLCore, EncrypterCore and FileHandler are the untrusted counterparts that doesn't reside in any enclave. These mainly consists of UI methods and file handling methods. The MLCore and EncrypterCore classes contains interfaces like secure login and gets inputs from the user, like the dataset uploader. These classes interact with the filehandler in the untrusted code base.

#### 4.4.1 Crypto Module

The crypto module of the application contains both trusted and untrusted code bases. This module requires the user to Login, access windows APIs and also provide enclave functions like data sealing, KDF etc. The crypto module resides completely on the local computer of the user with Intel SGX support. Figure 29 shows the communication and interaction between various components inside this module.

A registered user is required to authenticate himself to the system using Windows SSPI and hence initiates a session. The user's hash is retrieved by the crypto module to be further used to generate a UserKey. This UserKey is used to encrypt all the FileKeys under this user. An user is allowed to encrypt as many files as he wishes. Each of these FileKeys have to be persisted along with the encrypted message to be transferred through an insecure medium. Hence the FileKeys are encrypted with the UserKeys for stronger authentication and accountability.

Once a user is authenticated, the interface requests for files to be encrypted. The user inputs all the files to be encrypted. The files are then encrypted using AES-GCM method provided by Intel SGX. AES-GCM is an authenticated encryption method which requires a mandatory authentication tag that is used in the decryption process. So an encrypted file comprises of the cipher text, the encrypted FileKey and the authentication tag. The KDF requires a random salt as input. The

---

<sup>7</sup><https://software.intel.com/content/www/us/en/develop/articles/intel-digital-random-number-generator-drng-software-implementation-guide.html>

crypto module uses the random generator method of the DRNG class to generate the necessary salt for the KDF. Since the crypto module resides on the local computer of the user, the data sealing is done using the Enclave Author which is defined with the help of the user hash.

The crypto module allows the user to save the encrypted file (.CT format) anywhere on the computer. The user can open the file with any text editor. But any change to the cipher text will be detected by the authenticated encryption algorithm and will result in decryption failure. Any replay attacks are detected using this system.

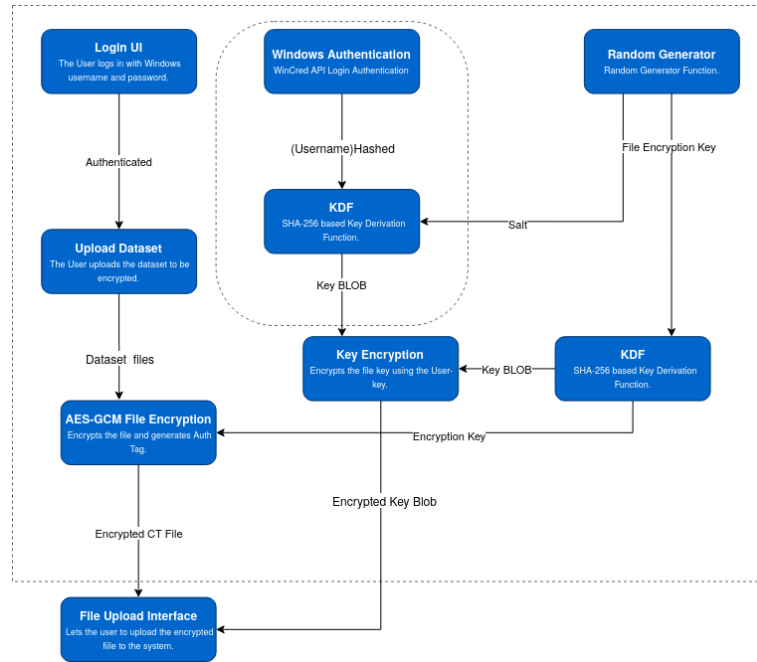


Figure 29: C2 - Crypto Module

On the other end, the decryption method parses the .CT file for cipher text, the encrypted FileKey and the authentication tag. The encrypted FileKey is first decrypted using the UserKey. The cipher text is then authenticated for its integrity and then decrypted by the AES decryption algorithm.

The plain text data after decryption doesn't leave the enclave and is hence secure from any attacks from malicious application in the cloud or the VMM, OS etc.

#### 4.4.2 Machine Learning Module

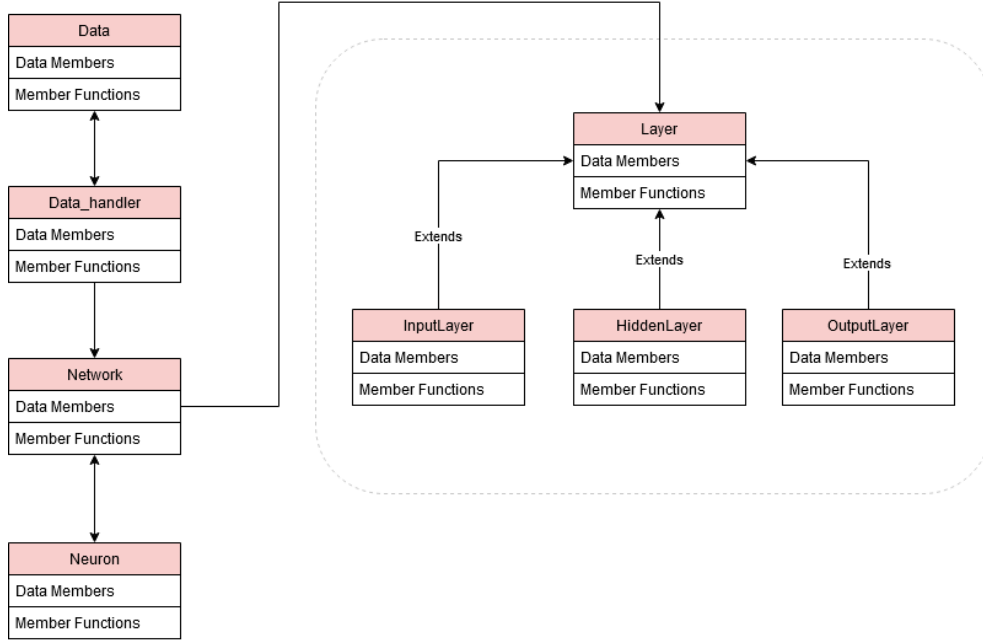


Figure 30: Machine Learning Enclave Code - Class Diagram

The ML module comprises of classes that handle the plain text data, classes for different layers and classes for individual neurons and networks of neurons. The ML architecture is a classic neural network setup with 3 different layers: input, hidden and output layer. The ML module requires the data to be in plain text as opposed to homomorphic schemes. The plain text data resides only in the enclave and is protected by attacks from the untrusted part of the cloud. No instance of the plain text data is passed outside the enclave in any form.

The setup has 50 hidden layers and trains for 1000 epochs. The algorithm reads in CSV files.

#### 4.5 Enclave Design

The enclave code is the most important part of the application design as it is the trusted part of the application. The goal is to minimize the number of components inside the enclave to reduce the

attack surface. In this design, there are 3 main components inside the enclave.

In the user end, the enclave is needed to seal the data safely. Hence two classes reside here. The Crypto class and the DRNG class. The enclave in the remote server (cloud) has all 3 components namely, Crypto, DRNG and machine learning classes (Section 4.4.2). These classes together make the enclave code so it is important that they are implemented effectively and are free of bugs.

#### 4.5.1 Enclave Logistics

The first step in designing the enclave is working out a system to manage the enclave itself. The enclave must be launched properly and the resulting enclave ID must be provided to the ECalls. This setup is ideal if this is transparent to the upper layers of the application. One easy solution is to use global variables in the Enclave Bridge DLL to hold enclave information. This design decision comes with a restriction: unless there is a locking mechanism of some sort, only one thread can be active in the enclave at a time. But for simple applications single pipelined threads are more than enough.

What happens if an enclave is lost due to a power event or a bug that causes it to crash? For that, we check the return value of the ECall: it indicates the success or failure of the ECall operation itself, not of the function being called in the enclave [45]. Any result other than `SGX_SUCCESS` indicates that the program did not successfully enter the enclave and the requested function did not run.

#### 4.5.2 Enclave Definition Language

When building an enclave using Intel SGX the interface has to be defined to the enclave in the EDL. An enclave's bridge functions, both its ECalls and OCalls, are prototyped in its EDL file. When the project is built Intel SGX SDK parses the EDL file and generates a series of proxy functions. Each ECall and OCall gets a pair of proxy functions: a trusted half and an untrusted half. The trusted functions go into `EnclaveProject.t.h` and `EnclaveProject.t.c` and are included in the Autogenerated Files folder of your enclave project. The untrusted proxies go into `EnclaveProject_u.h`

and EnclaveProject\_u.c and are placed in the Autogenerated Files folder of the project that will be interfacing with your enclave. The proxy functions are responsible for

1. Marshaling data into and out of the enclave
2. Placing the return value of the real ECall or OCall in an address referenced by a pointer parameter
3. Returning the success or failure of the ECall or OCall itself as an `sgx_status_t` value (Success or error).

Its general structure of an EDL file is as follows:

---

```
1  enclave{
2      // Include files
3      // Import other edl files
4      // Data structure declarations to be used as parameters
5      // of the function prototypes in edl
6
7      trusted {
8          // Include file if any. It will be inserted in the
9          // trusted header file (enclave_t.h)
10         // Trusted function prototypes (ECalls)
11     };
12
13     untrusted {
14         // Include file if any. It will be inserted in the
15         // untrusted header file (enclave_u.h)
16         // Untrusted function prototypes (OCalls)
17     };
```

---

## 5 Security Considerations

TEEs in general are not a new technology. Hundreds of millions of mobile devices worldwide rely on TEEs for protection of security-critical applications like DRM and OS components like Android Keystore. A TEE that is widely used in android devices is the ARM TrustZone. Unlike the Intel SGX, the ARM TrustZone uses a different architecture and implementation ideologies. But at the same time, TEEs have been successfully attacked with highly damaging impacts across various platforms. Unfortunately, these attacks have been possible not just by the presence of security flaws in TEE systems, but also the flaws and bugs in the application that runs in them. This section focuses on the security aspects of Intel SGX. It is to be noted that most of these points are applicable for other TEEs as well.

This section explains the security aspects, both positive and negative about this work. The second section explains the external authentication system(windows SSPI) used in the application and its associated security considerations. Finally, the possibility of side channel attacks in which a malicious user can get in to the system or gain access to the application security is explored.

### 5.1 Threat Model

For this work we assume the standard SGX threat model. The adversary is in complete control of privileged software, in particular, the hypervisor and the operating system. They are capable to spawn and stop processes at any point, set their affinity and modify it at runtime. The adversary can also read and write to any memory region except the enclave memory, map any virtual page to any physical one and dynamically re-map pages. The attacker also managed to have control of the cloud and any communication from the participating users and the cloud.

## 5.2 Attack surfaces

To preface the security aspects of this work, it is important to know how Intel Software Guard Extensions Technology helps secure data.

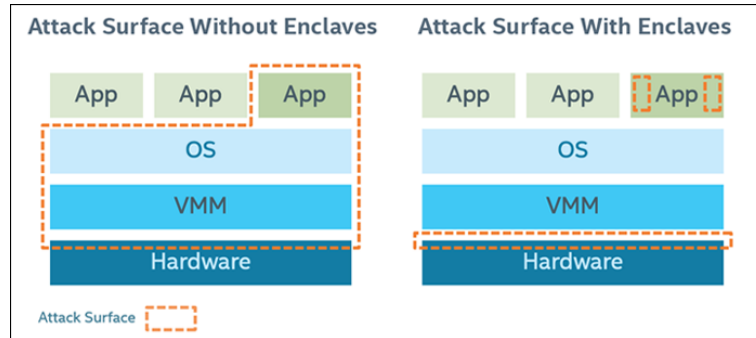


Figure 31: Attack surface with and without Intel SGX.

The Figure 31 shows the attack surface for a traditional application and for a enclave application.

Intel SGX offers the following protections from known hardware and software attacks:

1. Reading and writing to an enclave memory cannot be done outside the enclave.
2. Many exploits use debugger vulnerabilities or misplace debug code to break an application. Intel SGX can only be debugged by the Intel SGX Debugger. Production enclaves cannot be debugged by any software or hardware debuggers.
3. Since enclave need Ecalls and Ocalls to enter and leave the enclave, classic function calls, jumps, register manipulations or stack manipulations will not work in an enclave. The only way to call enclave functions is through a new instruction which performs several protection checks.
4. Enclave memory is encrypted using industry standard authenticated encryption which protects against replay attacks. This also means that trying to tap on to the memory directly or using the memory chip in another CPU will only yield encrypted data.

5. The memory encryption key is randomly changed between every power cycle. This key is then stored within the CPU and is not accessible by any other software.
6. Data stored within enclaves can only be accessed by code that shares the enclave.

In addition to the intrinsic security limitations of SGX documented by Intel, attackers can easily compromise the security of an SGX-based application by exploiting the bugs in the enclave software such as classical memory corruption or concurrency issues (SGX will not magically secure insecure software). Unsafe SGX programming by the enclave writer, for example by mishandling in-enclave pointers, or by leaking secret data through function calls outside the enclave.

Focusing on the development process, an insecure development environment can make the enclave application vulnerable to exploits. If the OS is compromised when developing the enclave then integrity of that enclave and any code within it is also compromised and should not be trusted.

Contrary to the suggestion by Intel to use 3rd party libraries, this can be a very serious vulnerability in the application. If the 3rd party library is compromised then so is the enclave code. Bugs in the Intel trusted libraries (libc or crypto library, for example) are also potential issues but these are addressed regularly by Intel.

Even though the CPU code cannot interact directly with the enclave, CPU bugs, in the hardware logic or in the microcode can be used to attack the enclave.

### 5.3 Windows SSPI

The Windows SSPI is a common interface between transport-level applications, such as Microsoft Remote Procedure Call (RPC), and security providers, such as Windows Distributed Security. SSPI allows a transport application to call one of several security providers to obtain an authenticated connection. In this work, the Windows SSPI is called by the win32 application to authenticate the user and initiate a secure session. Security packages support security protocols such as Kerberos authentication and the Microsoft LAN Manager.



SSPI allows an application to use various security models available on a computer or network without changing the interface to the security system. The SSPI doesn't create or handle logon credentials. This function is a privileged operation and is handled by the OS.

There are various authentication packages<sup>8</sup> provided by Microsoft. Some packages are the following.

1. Credential Security Support Provider
2. Microsoft Negotiate
3. Microsoft NTLM
4. Microsoft Kerberos
5. Microsoft Digest SSP
6. Secure Channel

For this work the Microsoft Credential SSP or CredSSP is used. CredSSP lets an application delegate the user's credentials from the client to the target server for remote authentication. CredSSP provides an encrypted Transport Layer Security Protocol channel. The client is authenticated over the encrypted channel by using the Simple and Protected Negotiate (SPNEGO) protocol with either Microsoft Kerberos or Microsoft NTLM. After the client and server are authenticated, the client passes the user's credentials to the server. The credentials are doubly encrypted under the SPNEGO and TLS session keys. CredSSP supports password-based logon as well as smart card logon.

Coupled with the Windows SSPI, the Winlogon module is used for an interactive logon. The Winlogon module uses windows 10 credential providers as the primary mechanism for user authentication. The credential providers are the only method for users to prove their identity which is required for logon and other system authentication services. There are several ways to logon to the application such as passwords, PIN, smartcard, fingerprint, Face and Iris recognition. This API

---

<sup>8</sup><https://docs.microsoft.com/en-us/windows/win32/secauthn/ssp-packages-provided-by-microsoft>

also includes support for 3rd party credential providers. The former are called System credential providers and the latter are called third-party credential providers.

To understand the security considerations of the authentication system as a whole, it is important to understand how the Windows credential management works. It is a process where the OS receives the credentials from the application, service or user and secures that information for future presentation to the authentication target. The credentials used in authentication are digitally verifiable documents that associate the user's identity to some form of proof of authenticity, such as a certificate, a password, or a PIN.

By default the Windows credentials are stored on the local computer and are validated against the Security Accounts Manager(SAM) database on the local computer as well. Credentials are requested and collected through user input on the logon user interface or coded through the application programming interface (API) to be presented to the authenticating target.

The credentials, for example the password is stored as an hash in the local system. Using the Windows credential manager increases the attack surface in the local system significantly. An attacker can pull credentials from many different areas of the system. Once the attacker gains access to a regular endpoint computer, the attacker can look for the credentials in

1. WDigest
2. SAM
3. Local Security Authority (LSA) secrets
4. Kerberos System

Even though the credentials stored here are encrypted or hashed, the attacker can attempt to pass it to the authentication server. If the server simply compares the encrypted credentials(hashed passwords), then it will match and the attacker will gain access. This vulnerability called "Pass the Hash"[48], where the attacker exploits the password hash which is saved in the cache memory or by stealing(credential dumping) in from the above listed places. This exploit lets the attacker

gain administrator-level permissions for a user's machine. On the other hand Windows credential manager is secure as long as the user themselves are to be trusted.

There are many ways to mitigate and defend against credential dumping. If databases like SAM, LSA are monitored and logged or keep a look out for unexpected connections from IP addresses not assigned to known domain controllers, reduces the risk of malicious activities.

What happens if an attacker already has all the credentials or if the endpoint computer cannot be trusted? How to defend the integrity and confidentiality of the system? The solution to this is Multi-Factor Authentication(MFA). MFA can serve as a second source of authentication system to verify user identity before granting access. Another way to defend against an attacker from moving laterally through a network and gain access to a machine is to implement a zero-trust strategy.

## 5.4 Side Channel Attacks

Side Channel attacks are those that target the vulnerability of a computer system and gain information from its implementation rather than the weakness of the software or algorithm. Timing information, power consumption, electromagnetic leaks or even sound can provide an extra source of information, which can be exploited. Since Intel SGX is a secure hardware based TEE, it is important to analyze potential side channel exploits that might compromise the confidential computing system.

The latest side channel exploits and research works are explored in Section 2.2.4. Common side channel attacks that might exploit enclaves are cache attacks, timing attacks and software-initiated fault attacks. Even though software initiated attacks are a class of side channel attacks, the work done in section 2.2.4 classifies software attacks as not a part of side-channel attacks. Spectre attacks [36] come under software initiated attacks which exploit speculative execution of enclave code. The cache timing attacks on enclaves by [39] shows how a malware can run inside an enclave as enclave code and compromise other enclaves in the same hardware by mounting a DRAM attack.

## 6 Conclusion

The rise of cloud technologies, security and privacy are the most significant concerns for both enterprises and consumers. The need for secure data and its computation is the primary motivation of this thesis. This work aims to minimize the data exposure to the intended authorized cloud subsystem only while maintaining integrity and confidentiality of both the data and its computation. Confidential computing aims to set a security standard to protect the data-in-use. It furthermore works on providing better user authority and transparency that is especially fitting to public clouds.

The purpose of this work is to propose a system of confidential computing using Intel SGX as a trusted environment and to study the practical implications of such a system. The work explores the design aspects of the system and also studies the security considerations of each of the components.

The state of the art in Chapter 2 reviews the technologies in the field of neural network based machine learning, confidential machine learning and concludes with TEEs and attacks on TEEs. The TEE chosen for this work is Intel's SGX technology which is the latest and more secure among the ones in market. A thorough background study on Intel SGX is in Chapter 3. Chapter 4 comprises of the implementation details of the practical work including the solutions offered to use cases and the entire software design of both trusted and untrusted parts of the system. The threat model and the security considerations of the system components are explained in Chapter 5.

Like any technology, Intel SGX is vulnerable to side channel attacks. Intel claims to have complete security from such attacks but many researches show many chinks in the SGX technology. Enclave codes are definitely safer than traditional applications and since it is a trust based system, Intel over the years has definitely earned this trust. It is important not to include numerous 3rd party libraries and plugins which has an adverse effect on the whole attacks surface of the system. Once the trust is partitioned and delegated to many components, the probability of exploits increases drastically. Responsible software design and development is key for minimizing risks system penetration and unauthorised access.

## References

- [1] MITRE Corporation, Common vulnerabilities and exposures (CVE) details: The ultimate security vulnerability datasource. browse vulnerabilities by date, 2019. Accessed on 10/10/2019. [Online]. Available: <https://www.cvedetails.com/browse-by-date.php>.
- [2] Rafael Pereira Pires, "Distributed systems and trusted execution environments: Trade-offs and challenges", arXiv:2001.09670v3 [cs.CR] 11 Mar 2020, Available: <https://arxiv.org/pdf/2001.09670.pdf>
- [3] D. Snyder, Intel unleashes next-gen enthusiast desktop pc platform at gamescom, 2015. Available: <https://blogs.intel.com/technology/2015/08/6th-gen-gamescom/>.
- [4] Redhat Corporation, Current Trusted Execution Environment landscape, December 2, 2019, Accessed on 07/05/2020. [Online]. Available: <https://next.redhat.com/2019/12/02/current-trusted-execution-environment-landscape/>
- [5] GENTRY, C. Fully homomorphic encryption using ideal lattices. In ACM Symposium on Theory of Computing (STOC) (2009).
- [6] Thore Graepel, Kristin Lauter, and Michael Naehrig<sup>1</sup>, ML Confidential: Machine Learning on Encrypted Data, Microsoft Research.
- [7] Zvika Brakerski and Vinod Vaikuntanathan. Fully homomorphic encryption from Ring-LWE and security for key dependent messages. In Phillip Rogaway, editor, CRYPTO, volume 6841 of Lecture Notes in Computer Science, pages 505–524. Springer, 2011.
- [8] Zvika Brakerski, Craig Gentry, and Vinod Vaikuntanathan. (Leveled) fully homomorphic encryption without bootstrapping. In Shafi Goldwasser, editor, Innovations in Theoretical Computer Science – ITCS 2012, pages 309–325. ACM, 2012.

- [9] Nigel P. Smart and Frederik Vercauteren. Fully homomorphic encryption with relatively small key and ciphertext sizes. In Phong Q. Nguyen and David Pointcheval, editors, Public Key Cryptography, volume 6056 of Lecture Notes in Computer Science, pages 420–443. Springer, 2010.
- [10] Sanjiv Ranjan Das, Data Science Theories, Models, Algorithms and Analytics. Available: <http://algo.scu.edu/sanjivdas/>
- [11] Kevin Gurney, An introduction to neural network, University of Sheffield, ISBN: 1-85728-673-1, UCL Press, 1997.
- [12] Andy Thomas, An Introduction to Neural Networks for Beginners. Accessed on : 12.05.2020 Available : <https://adventuresinmachinelearning.com/wp-content/uploads/2017/07/An-introduction-to-neural-networks-for-beginners.pdf>
- [13] Saad ALBAWI , Tareq Abed MOHAMMED, Saad AL-ZAWI, Understanding of a Convolutional Neural Network, ICET2017, Antalya, Turkey
- [14] TCG: Trusted computing group, Trusted platform modules (TPM) part 1, design principles, 2011. [Online]. Available: [https://trustedcomputinggroup.org/wp-content/uploads/TPM-Main-Part-1-Design-Principles\\_v1.2\\_rev116.01032011.pdf](https://trustedcomputinggroup.org/wp-content/uploads/TPM-Main-Part-1-Design-Principles_v1.2_rev116.01032011.pdf).
- [15] ARM security technology - building a secure system using TrustZone technology, ARM Limited, 2009. [Online]. Available: [http://infocenter.arm.com/help/topic/com.arm.doc.pr029-genc-009492c/PRD29-GENC-009492\\_C\\_trustzone\\_security\\_whitepaper.pdf](http://infocenter.arm.com/help/topic/com.arm.doc.pr029-genc-009492c/PRD29-GENC-009492_C_trustzone_security_whitepaper.pdf).
- [16] S. Pinto and N. Santos, “Demystifying arm TrustZone: A comprehensive survey”, ACM Computing Surveys, vol. 51, no. 6, 130:1–130:36, 2019. doi: 10.1145/3291047
- [17] AMD, EXTENDING SECURE ENCRYPTED VIRTUALIZATION WITH SEV-ES, KVM FORUM -2018 Available : <https://events19.linuxfoundation.org/wp-content/uploads/2017/12/Extending-Secure-Encrypted-Virtualization-with-SEV-ES-Thomas-Lendacky-AMD.pdf>

- [18] AMD, Developer Central, Available : <https://developer.amd.com/sev/>
- [19] David Kaplan, Jeremy Powell, Tom Woller, AMD Memory Encryption. Available : [https://developer.amd.com/wordpress/media/2013/12/AMD\\_Memory\\_Encryption\\_Whitepaper\\_v7-Public.pdf](https://developer.amd.com/wordpress/media/2013/12/AMD_Memory_Encryption_Whitepaper_v7-Public.pdf)
- [20] Yoongu Kim, Ross Daly, Jeremie Kim, Chris Fallin, Ji Hye Lee, Donghyuk Lee, Chris Wilkerson, Konrad Lai, Onur Mutlu, Carnegie Mellon University and Intel Labs, Flipping Bits in Memory Without Accessing Them: An Experimental Study of DRAM Disturbance Errors. Available: [https://people.inf.ethz.ch/omutlu/pub/dram-row-hammer\\_isca14.pdf](https://people.inf.ethz.ch/omutlu/pub/dram-row-hammer_isca14.pdf)
- [21] intel Developer Zone, Accessed on 15.05.2020, Available : <https://software.intel.com/en-us/forums/intel-software-guard-extensions-intel-sgx/topic/737218>
- [22] MERKEL, D. Docker: Light weight Linux Containers for Consistent Development and Deployment. Linux Journal (Mar.2014).
- [23] KUBERNETES. Available : <http://kubernetes.io>, 2016
- [24] Sergei Arnautov, Bohdan Trach, Franz Gregor, Thomas Knauth, Andre Martin, Christian Priebe, Joshua Lind, Divya Muthukumaran, Dan O’Keeffe, Mark L Stillwell, David Goltzsche, David Eysers, Rudiger Kapitza, Peter Pietzuch, and Christof Fetzer, SCONE: Secure Linux Containers with Intel SGX.
- [25] M.Abadi, P.Barham, J.Chen, Z.Chen, A.Davis, J.Dean, M.Devlin, S.Ghemawat, G.Irving, M.Isard, et al. Tensorflow: A system for large-scale machine learning. In OSDI, volume 16, pages 265–283, 2016.
- [26] Roland Kunkel, Do Le Quoc, Franz Gregor, Sergei Arnautov, Pramod Bhatotia, Christof Fetzer, TensorSCONE: A Secure Tensor Flow Framework using Intel SGX, arXiv:1902.04413v1 [cs.CR] 12 Feb 2019.

- [27] Zhongshu Gu et al., YerbaBuena: Securing Deep Learning Inference Data via Enclave-based Ternary Model Partitioning, IBM Research, arXiv:1807.00969v2 [cs.CR] 16 May 2019
- [28] H. Shacham, “The geometry of innocent flesh on the bone: Return-into-libc without function calls (on the x86)”, in Proceedings of the 14th ACM Conference on Computer and Communications Security, ser. CCS ’07, Alexandria, Virginia, USA, 2007, pp. 552–561. doi: 10.1145/1315245.1315313
- [29] T. Bletsch, X. Jiang, V. W. Freeh and Z. Liang, “Jump-oriented programming: A new class of code-reuse attack”, in Proceedings of the 6th ACM Symposium on Information, Computer and Communications Security, ser. ASIACCS ’11, Hong Kong, China: ACM, 2011, pp. 30–40. doi: 10.1145/1966913.1966919. [Online]. Available: <http://doi.acm.org/10.1145/1966913.1966919>
- [30] A. Baumann, M. Peinado and G. 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, 2014, pp. 267–283. [Online]. Available: <https://www.usenix.org/conference/osdi14/technicalsessions/presentation/baumann>
- [31] C. che Tsai, D. E. Porter and M. Vij, “Graphene-SGX: A practical library OS for unmodified applications on SGX”, in 2017 USENIX Annual Technical Conference (USENIX ATC 17), Santa Clara, CA: USENIX Association, 2017, pp. 645–658. [Online]. Available: <https://www.usenix.org/conference/atc17/technical-sessions/presentation/tsai>.
- [32] Daniel Moghimi, Berk Sunar, Thomas Eisenbarth, and Nadia Heninger, TPM-FAIL:TPM meets Timing and Lattice Attacks, 29th USENIX Security Symposium (USENIX Security 20), Available: <https://www.usenix.org/conference/usenixsecurity20/presentation/moghimi>
- [33] Eyal Ronen, Robert Gillham, Daniel Genkin, Adi Shamir, David Wong, and Yuval Yarom. The 9 Lives of Bleichenbacher’s CAT: New Cache ATtacks on TLS Implementations. In IEEE Symposium on Security and Privacy, 2019.



- [34] Kit Murdock, David Oswald, Flavio D. Garcia, Jo Van Bulck, Daniel Gruss, and Frank Piessens, Plundervolt: Software-based Fault Injection Attacks against Intel SGX, Proceedings of the 41st IEEE Symposium on Security and Privacy (S&P'20), Available : <https://plundervolt.com/>
- [35] Jo Van Bulck, Marina Minkin, Ofir Weisse, Daniel Genkin, Baris Kasikci, Frank Piessens, Mark Silberstein, Thomas F. Wenisch, Yuval Yarom, and Raoul Strackx, Foreshadow - Extracting Keys to the Intel SGX Kingdom with Transient Out-of-Order Execution, 27th USENIX Security Symposium. August 15–17, 2018 • Baltimore, MD, USA ISBN 978-1-939133-04-5
- [36] KOCHER, P., GENKIN, D., GRUSS, D., HAAS, W., HAMBURG, M., LIPP, M., MANGARD, S., PRESCHER, T., SCHWARZ, M., AND YAROM, Y. Spectre attacks: Exploiting speculative execution. arXiv preprint arXiv:1801.01203 (2018).
- [37] Guoxing Chen, Sanchuan Chen, Yuan Xiao, Yinqian Zhang, Zhiqiang Lin, Ten H. Lai, SgxPectreAttacks: Stealing Intel Secrets from SGX Enclaves via Speculative Execution, arXiv:1802.09085v3 [cs.CR] 3 Jun 2018
- [38] Johannes Götzfried et al. Cache Attacks on Intel SGX, EuroSec'17: Proceedings of the 10th European Workshop on Systems Security, April 2017, Available: <https://doi.org/10.1145/3065913.3065915>
- [39] Michael Schwarz, Samuel Weiser, Daniel Gruss, Clémentine Maurice and Stefan Mangard, Malware Guard Extension: abusing Intel SGX to conceal cache attacks. Available : <https://doi.org/10.1186/s42400-019-0042-y>
- [40] Kocher PC, Timing Attacks on Implementations of Diffe-Hellman, RSA, DSS, and Other Systems. In: Crypto'96. Available: [https://doi.org/10.1007/3-54068697-5\\_9](https://doi.org/10.1007/3-54068697-5_9)
- [41] D. Page, Theoretical Use of Cache Memory as a Cryptanalytic Side-Channel, University of Bristol Technical Report CSTR-02-003, Submitted to TISSEC, Available: <https://eprint.iacr.org/2002/169>

- [42] Bernstein DJ, Cache-timing attacks on AES. Technical report. Department of Mathematics, Statistics, and Computer Science, University of Illinois, Chicago. DOI: 10.1007/978-3-662-44709-3\_5
- [43] COLIN PERCIVAL, CACHE MISSING FOR FUN AND PROFIT, Available: [https://papers.freebsd.org/2005/cperciva-cache\\_missing/](https://papers.freebsd.org/2005/cperciva-cache_missing/)
- [44] Intel, Intel® Software Guard Extensions, Developer Guide, Available: [https://01.org/sites/default/files/documentation/intel\\_sgx\\_sdk\\_developer\\_reference\\_for\\_linux\\_os.pdf.pdf](https://01.org/sites/default/files/documentation/intel_sgx_sdk_developer_reference_for_linux_os.pdf.pdf)
- [45] Victor Costan and Srinivas Devadas, Intel SGX Explained, Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology.
- [46] NIST, Recommendation for Block Cipher Modes of Operation: Galois/Counter Mode (GCM) and GMAC, NIST Special Publication 800-38D November, 2007
- [47] D. Eastlake and P. Jones. RFC 3174: US Secure Hash Algorithm 1 (SHA1). Internet RFCs, 2001.
- [48] <https://www.blackhat.com/docs/us-15/materials/us-15-Moore-Defeating%20Pass-the-Hash-Separation-Of-Powers-wp.pdf>