

Plinius: Secure and Persistent Machine Learning Model Training

Peterson Yuhala
University of Neuchâtel
Neuchâtel, Switzerland
peterson.yuhala@unine.ch

Pascal Felber
University of Neuchâtel
Neuchâtel, Switzerland
pascal.felber@unine.ch

Valerio Schiavoni
University of Neuchâtel
Neuchâtel, Switzerland
valerio.schiavoni@unine.ch

Alain Tchana
ENS Lyon, France
Inria
alain.tchana@ens-lyon.fr

Abstract—With the increasing popularity of cloud based machine learning (ML) techniques there comes a need for privacy and integrity guarantees for ML data. In addition, the significant scalability challenges faced by DRAM coupled with the high access-times of secondary storage represent a huge performance bottleneck for ML systems. While solutions exist to tackle the security aspect, performance remains an issue. Persistent memory (PM) is resilient to power loss (unlike DRAM), provides fast and fine-granular access to memory (unlike disk storage) and has latency and bandwidth close to DRAM (in the order of ns and GB/s, respectively). We present PLINIUS, a ML framework using Intel SGX enclaves for secure training of ML models and PM for fault tolerance guarantees. PLINIUS uses a novel mirroring mechanism to create and maintain (i) encrypted mirror copies of ML models on PM, and (ii) encrypted training data in byte-addressable PM, for near-instantaneous data recovery after a system failure. Compared to disk-based checkpointing systems, PLINIUS is $3.2\times$ and $3.7\times$ faster respectively for saving and restoring models on real PM hardware, achieving robust and secure ML model training in SGX enclaves.

I. INTRODUCTION

Privacy-preserving machine-learning [7] is a challenging computational paradigm and workflow. Data and computation must be protected from several threats, *e.g.*, powerful attackers, compromised hypervisors and operating systems, and even malicious cloud or human operators [33], [34]. Preserving the confidentiality of the models (*i.e.*, weights and biases) being trained, as well as the input datasets is paramount: these are the most valuable business assets. Example application domains include health, finance, Industry 4.0, *etc.* Given the significant amount of computing resources typically required during the training phase, moving the model training over public clouds appears to be a pragmatic approach. However, this immediately leads to contradictory arguments. On the one hand, one benefits from the endless scalability and dependability features of public clouds, as well as pushing away valuable assets from potentially compromised on-premises infrastructures to the cloud. On the other hand, exposing confidential datasets and models to untrusted clouds must be avoided. Fig. 1 shows our target scenario, *i.e.*, training of ML models over untrusted public clouds, showcasing the threats that this scenario implies.

Trusted execution environments (TEE) are quickly becoming the go-to solution to tackle such confidentiality requirements, and several cloud providers nowadays offer TEE-enabled

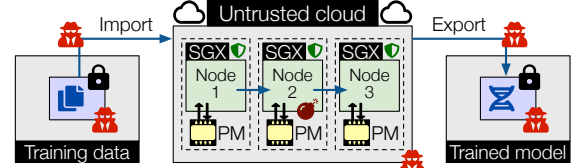


Fig. 1: ML and model building over untrusted public clouds. In case of failures or job pre-emption (●), the current state should be securely persisted to PM. Without proper mitigations, data and trained models could leak (⚠).

computing instances (IBM,¹ Azure²). Intel software guard extensions (SGX) [12] is a TEE that offers applications secure memory regions called *enclaves* to shield code and data from unwanted accesses. SGX is a promising candidate for protecting applications, given its wide availability across a variety of cloud providers. However, SGX imposes security restrictions on enclave code (*e.g.*, disallowing system calls), typically requiring application-level changes, as well as limited memory capacity of SGX enclaves, which requires developers to minimize the *trusted computing base* (TCB).

While offloading ML training jobs to SGX-enabled clouds might solve the confidentiality issue, such jobs are typically deployed in batch. Typically, batch jobs have lower priorities than latency sensitive (*e.g.*, production, user-facing) services [21]. In order to avoid resource waste because of workload variations, ML applications are colocated with latency sensitive applications. The former are automatically killed when the latter needs more resources. Another practice which may lead to interruptions of ML jobs is the use of cheap yet unreliable virtual machine instances such as EC2 spot instances [38]. The latter are automatically terminated when a better offer (*i.e.*, spot price) is made by another user. To avoid restarting model training from scratch when a task is killed, one can checkpoint/restore the current model on persistent storage. For instance, AWS SageMaker [28] suggests to frequently checkpoint the models to avoid abrupt data-loss when using spot instances. Fig. 3 represents this general ML pipeline. However, frequent checkpointing on secondary storage leads to significant I/O overheads, and relying on volatile memory (*i.e.*, DRAM) to mitigate these overheads would prevent to resume the job in case of task eviction.

¹<https://www.ibm.com/cloud/data-shield>

²<https://azure.microsoft.com/en-us/solutions/confidential-compute/>

Emerging memory technologies like *persistent memory* (PM) [40] have the potential to address the significant scalability challenges faced by DRAM, as well as the high latency of secondary storage. PM is persistent on power failure, byte-addressable, and can be accessed via processor `load` and `store` instructions. Recent work [40] shows how on-the-market PM solutions such as Intel Optane DC PM [40] result in significant performance gains for various applications. Cloud services like MS Azure already provide PM offerings [1], and we expect this technology to gain even more momentum. However, using PM in privacy-preserving ML jobs opens additional security risks: confidential model parameters could be persisted in plain-text on PM, or possibly be exposed at runtime to malicious privileged users or compromised operating systems. We take the stance that there is the need to develop tools and mechanisms to enable these applications to leverage PM in such secure computation environments. In this work we build the first framework that integrates secure ML with Intel SGX with fault tolerance on PM. State-of-the-art PM libraries (e.g., Intel Persistent Memory Development Kit⁵, Romulus [11], etc.), as well as ML frameworks (e.g., Tensorflow [6], Darknet [3], etc.), require considerable porting efforts to be fully functional within SGX enclaves. Tools exist (e.g., library OSes like Graphene-SGX [35] and Occlum [32], or modified enclave-compatible C libraries like SCONE [8]) to run unmodified applications inside SGX enclaves, at the downside of larger TCB sizes (thus larger attack surfaces) and large memory footprint, thus reducing performance. In ML scenarios with large confidential models and data sets, enclave memory becomes a major bottleneck that only important engineering efforts could optimize.

We present PLINIUS, a secure ML framework that leverages PM for fast checkpoint/restore of machine learning models. PLINIUS leverages Intel SGX to ensure confidentiality and integrity of ML models and data during training, and PM for fault tolerance. PLINIUS employs a *mirroring mechanism* which entails creating an encrypted mirror copy of an enclave model directly in PM. The mirror copy in PM is synchronized with the enclave model across training iterations. PLINIUS maintains training data in byte addressable PM. Upon a system crash or power failure while training, the encrypted ML model replica in PM is securely decrypted in the enclave, and used as next starting point of the training iteration: the training resumes where it left off, using training data already in memory. This avoids costly serialization operations of disk-based solutions. To validate our approach, we build and contribute SGX-DARKNET, a complete port of Darknet ML framework [3] to SGX, as well as SGX-ROMULUS, an SGX-compatible PM library on top of an efficient PM library [11]. We compare SGX-ROMULUS with unmodified Romulus library running in a SCONE container and our results show SGX-ROMULUS is best suited for PLINIUS framework. Using PLINIUS, we build and train *convolutional neural network* (CNN) models with real world datasets (i.e., MNIST [2]) and show PLINIUS reduces overhead by $\sim 3.5\times$ for model saving, and $\sim 2.5\times$ for model restores with real SGX

hardware and emulated PM.³

In summary, we make the following contributions:

- We implement and release as open-source SGX-ROMULUS⁴ on top of Romulus [11] for Intel SGX. SGX-ROMULUS manipulates PM directly from within SGX enclaves, without costly enclave transitions between secure and unsecured parts of an SGX application.
- We design, build, and release as open-source, SGX-DARKNET⁴ an extension of Darknet [3] for Intel SGX. SGX-DARKNET can perform secure training and inference on ML models directly inside SGX enclaves.
- We present PLINIUS, an open-source framework⁴ that leverages SGX-ROMULUS and SGX-DARKNET to provide an end-to-end fault tolerance mechanism to train models in privacy preserving ML settings.
- We provide a comprehensive evaluation of PLINIUS, using real PM hardware and real AWS Spot traces, showing its better performance when compared with traditional checkpointing on secondary storage (i.e., disk or SSD)

Roadmap. This paper is organized as follows. §II describes background concepts, while §III presents our threat model. §IV presents the architectures of SGX-ROMULUS, SGX-DARKNET and PLINIUS, while §V digs into the implementation details. The experimental evaluation of our system is in §VI. We discuss related work in §VII, before concluding and hinting at future work in §VIII.

II. BACKGROUND

This section presents a background on Intel SGX, PM, as well as some machine-learning concepts specific to PLINIUS.

Intel software guard extensions (SGX) [12] is a set of extensions to Intel’s architecture that permits applications to create CPU-protected memory areas (i.e., *enclaves*) shielding confidential code and data from disclosure and modifications.

SGX reserves a secure memory region called the *enclave page cache* (EPC) for enclave code and data. The processor ensures that software outside the enclave (e.g., the OS kernel or hypervisor) cannot access EPC memory. The enclave can access both EPC and non-EPC memory.

Data in the EPC is in plaintext only in on-chip caches and is encrypted and integrity-protected in the *memory encryption engine* (MME) once it is evicted from the cache to memory. Current Intel processors support a maximum of 128 MB of EPC memory, of which 93.5 MB is usable by SGX enclaves. This limits the total size of code and data allowed within the EPC. To support applications with larger memory needs, the Linux kernel provides a paging mechanism for swapping pages between the EPC and untrusted memory.

Enclaves cannot issue system calls and standard OS abstractions (e.g., file systems, network), which are ubiquitous in real world applications. All system services thus require costly enclave transitions, up to 13’100 CPU cycles [39]. The Intel SGX application design requires splitting applications

³Machines with SGX and PM support not on the market yet (Oct/2020).

⁴<https://github.com/Yuhala/sgx-pm-ml>

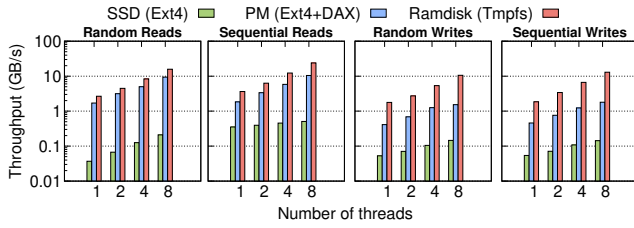


Fig. 2: Read/write throughput for sequential/random workloads on SSD, PM and Ramdisk using the `sync` I/O engine on FIO. 512 MB file per thread, 4 KB block size. Write workloads issue an `fsync` for each written block, average over 3 runs.

into a trusted (the enclave) and untrusted part. To achieve communication across the enclave boundary, the Intel SGX SDK provides specialized function call mechanisms, *i.e.*, `ecalls` and `ocalls`, respectively to enter and exit an enclave [10].

To mitigate security risks, the TCB should be as small as possible. Systems exist to run unmodified applications inside enclaves, either by porting entire library OSes into the enclave [32], [35] or via a modified `libc` library, specialized for containerized services [8]. These solutions are efficient with small application binaries, but quickly show limits for memory-constrained applications such as ML.

Persistent memory and PM libraries. PM is a novel memory technology that is *non-volatile*, *byte-addressable*, and has latency and bandwidth similar to that of DRAM. PM resides on the memory bus and can be accessed directly using CPU `load` and `store` instructions. Intel Optane DC PM is commercially available since April 2019. Optane DC PM scales better than DRAM and hence provides much larger capacity (up to 512 GB per PM module). Intel Optane DC PM modules can operate in two modes: *memory mode* where they are simply used to expand main memory capacity without persistence, and *app direct mode* which provides persistence [40]. PLINIUS leverages PM in *app direct mode*.

Applications can leverage PM in *app direct mode* by using standard operating system calls (*i.e.*, `read`, `write`) through the file system in the same way slower storage devices like SSDs and HDDs are accessed. This improves application performance but does not leverage the load/store interface provided by the PM modules. In PLINIUS, we enhance the ML library to do direct loads/stores from/to PM. This configuration is more challenging as it requires application modification. However, it results in more significant performance gains since persistent updates bypass both the kernel and file system [22].

Server-grade CPUs natively support up to 3 TB of PM [27], hence revealing PM as an attractive solution for fault tolerant applications. Due to data remanence [41], using PM could introduce security risks, in particular for confidentiality and data integrity.

The use of PM requires a paradigm shift for application developers. Several software tools and libraries have been proposed, such as Romulus [11], Mnemosyne [36] or Intel’s PMDK⁵ to facilitate PM related development. PM libraries expose PM to applications by memory-mapping files on a

persistent memory-aware file system with *direct access* (DAX) capabilities. DAX removes the OS page cache from the I/O path and allows for direct access to PM with byte-granularity.⁵

To characterize our PM units, we execute FIO⁶ with sequential and random workloads, and compare the read and write throughputs for native Ext4 over an SSD drive, Ext4+DAX on PM, and a `tmpfs` partition over volatile DRAM. We observe (Fig. 2) that the DAX-enabled file system on PM performs consistently better than its non-DAX counterpart on SSD, and is close to RAM-`tmpfs` performance (in the order of GB/s).

Specific processor instructions (*i.e.*, `CLFLUSH`, `CLFLUSHOPT`, `CLWB`) are used to flush data from cache lines to the PM memory controller. Through asynchronous DRAM refresh [40] data in the memory controller’s write buffers is guaranteed to be persisted in PM in case of a power failure. Persistence fences (*i.e.*, `SFENCE`) guarantee consistency by preventing `store` instructions from being re-ordered by the CPU. PM libraries like Romulus and the PMDK provide transactional API which enable developers to perform atomic updates on persistent data structures.

Romulus provides durable transactions via twin copies of data in PM and relies on a volatile log to track memory locations being modified in a transaction. The first copy, called the `main` region, is where user-code executes all in-place modifications; the second copy, the `back` region, is a *backup* (or snapshot) of the previous consistent state of the `main` region. Following a crash while mutating `main`, the content of `back` is restored to `main`. Romulus uses at most four persistence fences for atomic updates on data structures, regardless of transaction size, and a *store interposition* technique to ensure cache lines are correctly flushed to PM. The design of Romulus permits to have low *write amplification* [11] relative to other PM libraries, and hence we build on top of it as PM library, by porting it to be SGX-compatible.

Training ML models. A ML *model* can be described as a function that maps an input to a target output based on a set of parameters [19]. A linear regression model for example is a function $f(\mathbf{x}) = \mathbf{W}^T \mathbf{x} + \mathbf{b}$ where \mathbf{W} represents the model weights, \mathbf{b} the bias vector, and \mathbf{x} the input vector. The weights and biases are the learnable parameters of a model, and they determine the output of the latter for a given input.

The goal of model training is to obtain the set of learnable parameters that minimizes a *loss function* and maximizes the model’s accuracy on the training data. The loss function is a scalar function that quantifies the difference between the predicted value (for a given input data point) and the ground truth or real value [6]. During training, the learning algorithm iteratively feeds the model with batches of training data, calculates the loss, and updates the model parameters in such a way as to minimize the loss. A very popular learning algorithm used in ML for loss minimization is *stochastic gradient descent* (SGD) [29].

In this work, we rely on supervised learning, where a (costly) training phase builds a model out of labelled data, followed

⁵<https://pmem.io/>

⁶<http://freshmeat.sourceforge.net/projects/fio>

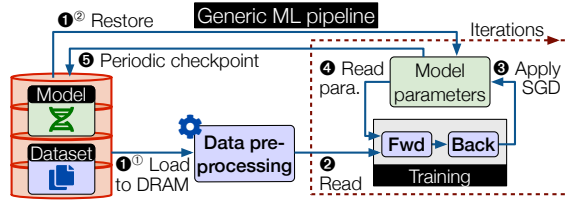


Fig. 3: General machine learning pipeline: *Fwd*=Forward propagation, *Back*=Backward propagation.

by a classification/inference phase using the model. Examples include visual object recognition, spam filtering, *etc.*

Fig. 3 shows a typical *supervised* ML model training pipeline. Training data is read from secondary storage (Fig. 3-1), preprocessed and used to train the model (Fig. 3-2,3,4). Training models such as deep neural networks can take up to several days, a time window sufficiently long for training jobs to experience failures or pre-emptions [6]. Also, large ML datasets (*i.e.*, order of GBs) are very common in the training phase. In the event of a failure during training, the model being trained as well as the training data sets resident in DRAM are lost and need to be re-read from secondary storage upon restart. Several state-of-the-art ML frameworks (*e.g.*, Tensorflow [6], Darknet [3], Caffe [24], *etc.*) provide mechanisms to checkpoint model states to secondary storage during training. However, the high latency and low bandwidth (in the order of MB/s) of secondary storage makes failure recovery a fundamental problem. The mentioned appealing properties of PM make the latter particularly interesting for fault tolerance in such ML scenarios.

In this work, we implement PLINIUS, a novel ML framework which leverages PM for fault tolerance and Intel SGX to ensure confidentiality and integrity of ML models, as well as sensitive training data. We build our ML framework on Darknet, which is popular in the ML community, provides good performance and is easily portable to Intel SGX.

III. THREAT MODEL

PLINIUS has three primary goals: (1) to ensure confidentiality and integrity of a ML model’s parameters (*e.g.*, weights, biases) during training; (2) to ensure confidentiality and integrity of the model’s replica on PM used for fault tolerance; and (3) to ensure confidentiality and integrity of training data in byte-addressable PM.

The system is designed to achieve these goals while facing a powerful adversary with physical access to the hardware and full control of the entire software stack including the OS and hypervisor. The adversary seeks sensitive information inside the enclave, on DRAM or PM, or data from the processor.

Model hyper-parameters such as model architecture, number of layers, size of training batches or type of training data are usually public information, as they do not leak any information about trained model parameters or sensitive training data [15], [17], [19]. In order to mitigate possible threats linked to malicious data sources, PLINIUS supports secure provisioning of model hyper-parameters via the SGX remote attestation mechanism. We assume that the adversary cannot physically

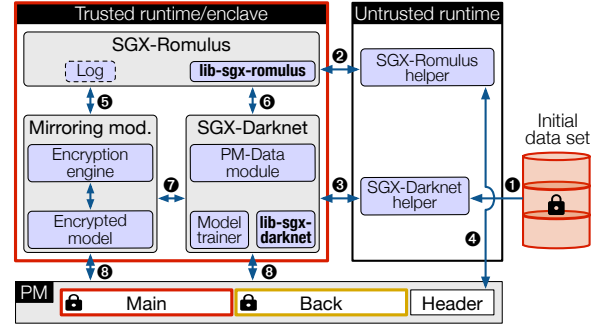


Fig. 4: PLINIUS architecture.

open and manipulate the processor package, that enclave code is correct and it does not leak sensitive information (*e.g.*, encryption keys) intentionally. Denial-of-service and side-channel attacks [9], [31], for which solutions exist [16], [30], are considered out of scope.

IV. PLINIUS ARCHITECTURE

The design of PLINIUS introduces an important issue: minimizing the TCB. A design approach based on a libOS like Graphene SGX or SCONE containers introduces thousands of lines of code into the enclave runtime, increasing security risks. Furthermore, with such a design, the enclave quickly reaches its memory limitation leading to a dramatic loss of performance. In light of these problems and following the SGX guidelines [10], we design an architecture partitioned into trusted and untrusted parts.

By manually porting the PM and ML libraries via separation into trusted and untrusted components, PLINIUS achieved a TCB reduction of $\sim 44\%$ in terms of LOC.

The architecture of PLINIUS consists of three main components interacting with each other: (1) an SGX-compatible PM library, *i.e.*, *sgx-romulus*; (2) an SGX-compatible deep-learning framework, *i.e.*, *sgx-darknet*; and (3) a mirroring module, which synchronizes the ML model inside the enclave with its encrypted mirror copy in PM. Figure 4 shows how these components interact. We detail each of them in the remainder of this section.

SGX-Romulus is a port of Romulus [11] to Intel SGX. SGX-ROMULUS implements durable transactions in PM directly within an SGX enclave. It consists of a secure user-space library, *lib-sgx-romulus*, which provides durable, concurrent transactions, and persistence primitives required to create and manage persistent data structures in PM. SGX-ROMULUS maintains a volatile log in enclave memory which logs the addresses and ranges of modified data in the current transaction. The size of the log varies with transaction size. A helper library in the untrusted runtime, *sgx-romulus-helper*, communicates with SGX-ROMULUS and permits to invoke necessary system calls (*e.g.*, *mmap*, *munmap*) which are required when leveraging PM via a DAX-enabled file system.

At application initialization, *sgx-romulus-helper* memory-maps the file corresponding to the persistent memory region (*main*, *back* and *header*) into application virtual address space (VAS), via a *mmap* system call.

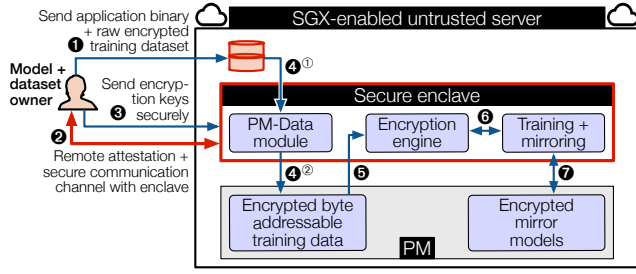


Fig. 5: Full model training workflow with PLINIUS.

Then, `sgx-romulus-helper` initializes the *persistent header* [11], which holds metadata to track the consistency state of the *main* and *back* regions, a reference to an array of persistent memory objects, and a pointer to the memory allocator’s metadata (e.g., allocated and unallocated PM). The address of the persistent header is passed to SGX-ROMULUS via an `ecall` and once the enclave validates this address, it then completes the PM region initialization, as further detailed in Algorithm 1. The enclave can then create or update persistent data structures in PM. Upon graceful termination of the enclave application, the enclave runtime issues an `ocall` to unmap the PM region from application VAS via the `munmap` system call.

SGX-Darknet is a port of Darknet [3] to Intel SGX. While many TEE-based ML libraries only provide support for inference, SGX-DARKNET supports both secure training and inference on ML models in Intel SGX enclaves. In order to achieve a minimal TCB, we separate SGX-DARKNET into trusted and untrusted parts. Our separation strategy involves keeping out of the enclave, as much as possible, computations which do not require any particular security. Examples include parsing of model configuration files, and initial data loading into DRAM.

To minimize code changes for commonly used (but unsupported) routines in Darknet (e.g., `fread`, `fwrite` etc.), SGX-DARKNET redefines the former as wrapper functions for `ocalls` to the corresponding libC functions in the untrusted runtime. A support library in the untrusted runtime, `sgx-darknet-helper`, provides the implementations of those `ocalls` invoking the corresponding libC routines. The *pm-data-module* is used by SGX-DARKNET to write/read encrypted data sets to/from PM. Finally, *lib-sgx-darknet* provides the API to train and do inference on models from within the enclave runtime.

Mirroring module. This component is in charge of creating and updating encrypted mirror copies of enclave ML models on PM. It contains the necessary logic to instantiate models that are both persistent and directly byte-accessible via loads and stores. It leverages the transactional API provided by SGX-ROMULUS to perform atomic updates on persistent models in PM. This is crucial as it prevents any inconsistency in PM data structures in the event of a system failure during data updates. The logic for building and managing persistent versions of complex data structures like ML models can get very bulky, and so we preferred to build a separate module for that rather than integrate it directly into SGX-DARKNET.

The *encryption engine* is responsible for encrypting/decrypting model parameters to be mirrored to/from the PM model, as well as in-enclave decryption of encrypted training data resident in PM.

In-enclave symmetric encryption/decryption relies on AES Galois counter mode (GCM) [13] implementation from the Intel SGX SDK. AES-GCM uses a 128, 192 or 256 bit key for all cryptographic operations, and provides assurance of the integrity of the confidential data [13]. PLINIUS uses a 128 bit key for all cryptographic operations.

As recommended by [13], for every encryption operation, we generate a random 12-byte *initialization vector* (IV) using the `sgx_read_rand()` [10, p. 200] function from the Intel SGX SDK. The encryption algorithm divides each plain text buffer into 128 bit blocks which are encrypted via AES-GCM. The IV and a 16-byte message authentication code (MAC) are then appended to each encrypted data buffer. The MAC is used to ensure data integrity during decryption.

The key used for encryption/decryption can be provisioned to the enclave via remote attestation [10, p. 99] or could be generated securely (e.g., if training data is not encrypted) inside the enclave using `sgx_read_rand()`. The encryption key, once generated or provisioned, can be securely sealed [10, p. 96] by the enclave for future use.

Full ML workflow with PLINIUS. Figure 5 shows the full ML workflow with PLINIUS. The owner of the data and the model sends the application binary and raw encrypted training data to the remote untrusted server (Figure 5-1). She then performs remote attestation (RA), establishes a secure communication channel (SC) with the enclave (Fig. 5-2) and sends encryption keys to the latter (Figure 5-3). The PM-data module transforms encrypted data on disk to encrypted byte addressable data in PM (Figure 5-4^{1,2}). The training module reads and decrypts (with keys obtained from RA & SC) batches of training data from PM (Figure 5, 5-5) with the trained model being mirrored to PM or into the enclave for restores (Figure 5-7).

Integration with different ML libraries. The current PLINIUS architecture uses Darknet as the ML library, due to its efficient and lightweight implementation in C that facilitates integration with SGX enclaves. Other ML libraries could be integrated into the PLINIUS architecture. In fact, once the ML library is ported to SGX, the same PLINIUS architecture holds.

To validate the generality of our architecture, we applied our mirroring mechanism within Tensorflow [6], another popular ML library. Tensorflow uses *tensor* data structures to store model information (e.g., weights and biases). Our implementation creates mirror copies of tensors in PM and restores them in enclave memory using PLINIUS’s mirroring mechanism. However, due to the large memory footprint of Tensorflow-based ML applications with respect to our EPC limit (93.5 MB), we opted to use Darknet ML library, which is lightweight but equally efficient.

Algorithm 1 — Initialization algorithms.

```
1 ## Untrusted (outside of enclave) ##
2 function init_sgx_romulus(pm_file)
3   mapped_addr = mmap(pm_file)
4   header_addr = create_header(mapped_addr)
5   ecall_init(header_addr)
6 end
7 function ocall_unmap
8   munmap(pm_file)
9 end
10 ## Trusted (inside enclave) ##
11 function ecall_init(header_addr)
12   initialize_main_and_back(header_addr)
13   recover() // defined in [11]
14 end function
```

V. IMPLEMENTATION DETAILS

We implement PLINIUS in C and C++ and it comprises 28'450 lines of code (LOC) in total, the trusted portion being 15'900 LOC. We use Intel SGX SDK v2.8 for Linux. The total size of application binary including the enclave shared library after compilation is 3 MB. In the remainder, we describe further details and a rundown for a ML model training.

Initialization. In this phase, PLINIUS memory maps PM into application virtual address space (VAS) (see Algorithm 1, lines 3-5) and initializes the persistent regions *main* and *back* (Algorithm 1, line 12), so that both regions are consistent before the training starts.

Initial dataset loading to PM. One key aspect of PLINIUS is the ability to use training data in PM. In PLINIUS, we load training data into PM once, after which the data stays in (byte addressable) PM. At resumption following a power failure or system crash, training data in PM is instantly accessible to the training algorithm, unlike in disk or SSD-based systems where data needs to be re-read from slow secondary storage into DRAM.

Initially, the training dataset is stored encrypted as files on secondary storage. Darknet training algorithms process input data as multidimensional arrays or matrices. The goal of this step is to load training data into such a data matrix in PM. The *sgx-darknet-helper* reads initial training data and labels from secondary storage into DRAM as a volatile matrix variable. The address of this matrix is sent to SGX-DARKNET via an *ecall*. The *pm-data-module* creates a corresponding persistent matrix on PM using the *lib-sgx-romulus* API. We annotate all persistent types (e.g., matrix rows, matrix values, model layer attributes, etc.) with the `persist<>` class from *lib-sgx-romulus*. This wrapper class ensures every `store` operation on the associated persistent data is followed by a *persistent write back* (PWB) to flush the cache line to PM. An appropriate fence instruction is used when ordering is required (e.g. at the end of a transaction).⁷ Once the persistent matrix is created, the training data is simply `memcpy`-ied from DRAM into PM within a transaction from within the enclave. The persistent data can then be accessed directly via its address.

Model training and mirroring. The PLINIUS architecture fits well for training neural network models [29], creating a

⁷Romulus supports 3 PWB + fence combinations: `clwb+sfence`, `clflushopt+sfence` (used in PLINIUS) and `clflush+nop`.

Algorithm 2 — Training a ML model in PLINIUS

```
1 function train_model(config)
2   enclave_model = create_enclave_model(config)
3   if not exists(pm_data) then
4     ocall_load_data_in_pm()
5   end if
6   iter = 0
7   if exists(pm_model) then
8     mirror_in(enclave_model)
9     iter = pm_model.iter
10  else
11    pm_model = alloc_mirror_model(enclave_model)
12  end if
13  while iter < MAX_ITER do // train for max_iter iterations
14    data_batch = decrypt_pm_data(batch_size)
15    train(enclave_model, data_batch)
16    mirror_out(enclave_model, iter)
17  end while
18  free(enclave_model)
19  ocall_unmap
20 end function
```

secure model in enclave memory. The architecture of the model and its hyper-parameters (e.g., layer types, batch size, learning rate, etc.) are defined in a config file which is parsed into a `config` data structure by *sgx-darknet-helper* in the untrusted runtime. Its address is sent to the enclave via an *ecall* where it is used to build the enclave model. If the training dataset has not been loaded in PM, an *ocall* is performed to load data once from secondary storage into a volatile data matrix variable accessible by the enclave runtime; this could be done in batches if the training dataset is very large. The training data is then loaded into PM (see §V).

If a persistent mirror model exists on PM, we *mirror-in* (read from PM and decrypt in enclave) its parameters into the enclave model, otherwise we allocate one in PM. Neural network models in general consist of multiple processing layers with learnable parameters (i.e., weights and biases). Darknet tracks these layer addresses in an array. In PLINIUS, we represent a neural network model on PM as a linked list of persistent layer structures, so as to simplify future modifications to the model's structure (e.g., add or remove layers). The model's layers contain persistent attributes, e.g., weight vector, bias vector, etc. These attributes are annotated with the `persist<>` class [11], which ensures PWBs are done for all stores to the corresponding persistent data.

Algorithms 2 and 3 summarize respectively model training and mirroring in PLINIUS. During model training, batches of training data are decrypted from PM (Algorithm 2, line 15) into enclave memory and used to train the enclave model for one training iteration. After each training iteration we do a *mirror-out* (encrypt in enclave and write to PM) of the enclave model parameters to its persistent mirror copy on PM. In the event of a crash during training, upon resumption the model and training data are already in PM and can be quickly `memcpy`-ied from PM into secure enclave memory. This obviates the need for much more slower reads from storage devices like SSDs and HDDs.

VI. EVALUATION

Our experimental evaluation of the PLINIUS prototype answers the following questions:

Algorithm 3 Mirroring algorithms.

```

1 function alloc_mirror_model(enclave_model)
2   BEGIN_TRANSACTION                                     // defined in [11]
3   head_pm_L = PMalloc(size)                             // L: neural network layer
4   head_pm_L.W = PMalloc(size)                           // W: layer's parameters
5   head_pm_L.next = nullptr
6   cur_pm_L = head_pm_L
7   n = enclave_model.numL
8   for i = 2 to n do
9     cur_pm_L.next = PMalloc(size)
10    cur_pm_L = cur_pm_L.next
11    cur_pm_L.W = PMalloc(size)
12    cur_pm_L.next = nullptr
13   end for
14   END_TRANSACTION                                     // defined in [11]
15 end function
16 function mirror_out(enclave_model, iter)
17   BEGIN_TRANSACTION
18   n = enclave_model.numL
19   pm_model.iter = iter
20   temp = head_pm_L
21   for i = 1 to n do
22     temp.W = encrypt(enclave_model.L(i).W)
23     temp = temp.next
24   end for
25   END_TRANSACTION
26 end function
27 function mirror_in(enclave_model)
28   n = pm_model.numL
29   temp = head_pm_L
30   for i = 1 to n do
31     enclave_model.L(i).W = decrypt(temp.W)
32     temp = temp.next
33   end for
34 end function

```

- How SGX-ROMULUS compares against unmodified Romulus in a Scone container ?
- How does PLINIUS improve checkpoint/restore performance when compared to secondary storage (e.g SSD)?
- How scalable is PLINIUS when varying model sizes?
- What are the main bottlenecks in the PLINIUS design?
- What is the overhead of batched-data decryptions?
- Is the mirroring mechanism robust against crashes?
- Are there processing and storage bottlenecks?

Experimental setup. At the time of this writing (October 2020), servers that support both SGX and PM are not available. Hence, we perform our experiments on two different servers, i.e., *sgx-emlPM* and *emlSGX-PM*. The *sgx-emlPM* node supports SGX but has no physical PM, hence we resort to emulating the latter with Ramdisk. This machine is equipped with a quad-core Intel Xeon E3-1270 CPU clocked at 3.80 GHz, and 64 GB of DRAM. The CPU ships with 32 KB L1i and L1d caches, 256 KB L2 cache and 8 MB L3 cache. Concerning *emlSGX-PM*, it is equipped with 4× Intel OptaneDC DIMMs of 128 GB each. However its processors lack native support for SGX. Hence, we resort to SGX in simulation mode [10]. The *emlSGX-PM* node is a dual-socket 40-core Intel Xeon Gold 5215 clocked at 2.50 GHz and 376 GB of DRAM. Each processor has 32 KB L1i and L1d caches, 1 MB L2 cache and a shared 13.75 MB L3 cache. Both servers run Ubuntu 18.04.1 LTS 64 bit and Linux kernel 4.15.0-54. We run the Intel SGX platform software, SDK and driver version v2.8. All our enclaves have max heap sizes of 8 GB and stack sizes of 8 MB. The EPC size is 128 MB (93.5 MB usable). Unless stated otherwise, we use CLFLUSHOPT and SFENCE for persistent write backs and ordering.

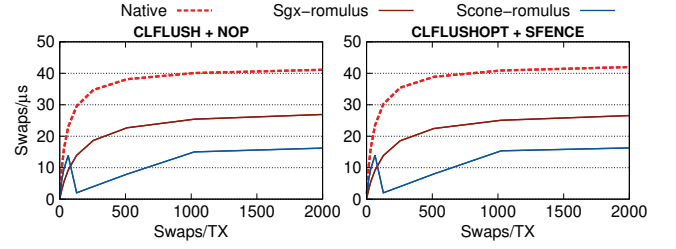


Fig. 6: SPS benchmark on the *sgx-emlPM* with varying transaction sizes for two PWB and fence combinations: CLFLUSH+ NOP (left) and CLFLUSHOPT + SFENCE (right).

By using both servers, we highlight the performance implications of both real SGX and real PM. All experimental comparisons are executed separately for each server, as they have completely different characteristics. We indicate where necessary on which node an experiment is carried out. All Scone containers are based on Alpine Linux [25].

All models used in our evaluations are convolutional neural networks (CNNs). The convolutional layers use *leaky rectified linear unit* (LReLU) [15] as activation, and all output layers are *softmax* [29] layers. The model optimization algorithm used is stochastic gradient descent (SGD), and the learning rate used is 0.1. Except stated otherwise, all training iterations use a batch size of 128. Concerning the dataset, we use MNIST [2], a very popular dataset in the deep learning community. It consists of 70'000 grayscale images of handwritten digits (60'000 training samples and 10'000 test samples).

Why SGX-Romulus makes sense. We begin by comparing SGX-ROMULUS with the unmodified Romulus library running in a Scone container, with the goal of understanding how a manually ported library using the Intel SGX SDK and the unmodified version in Scone stand against each other.

We measure how many swaps per second (SPS) they achieve, a metric commonly used [11] to compare PM libraries. SPS stores an array of integers in PM and evaluates the overhead of randomly swapping array values within a transaction, for different *persistence fences* and transaction sizes. This experiment uses the *sgx-emlPM* node, as real SGX is the main factor that dictates the performance differences. For each transaction size we run SPS for 20s. Figure 6 shows the throughput of swap operations on a 10 MB persistent array with different transaction sizes for different systems with a single threaded application. We include results for two choices of PWB implemented by Romulus and SGX-ROMULUS: *clflush + nop* and *clflushopt + sfence*. Our servers do not have support for *clwb*.

We observe that in both cases, the persistence fences take approximately 1.6× to 3.7× longer to complete in SGX-ROMULUS when compared to native (no SGX) systems for transaction sizes between 2 and 2048 swaps operations per transaction. When compared to Romulus in Scone, transactions for both fence implementations in SGX-ROMULUS are approximately 1.5× to 2.5× slower for transaction sizes between 2 to 64 swap operations per transaction. Beyond 64 swap operations per transaction, there is a pronounced drop in throughput for Romulus in Scone, and SGX-ROMULUS

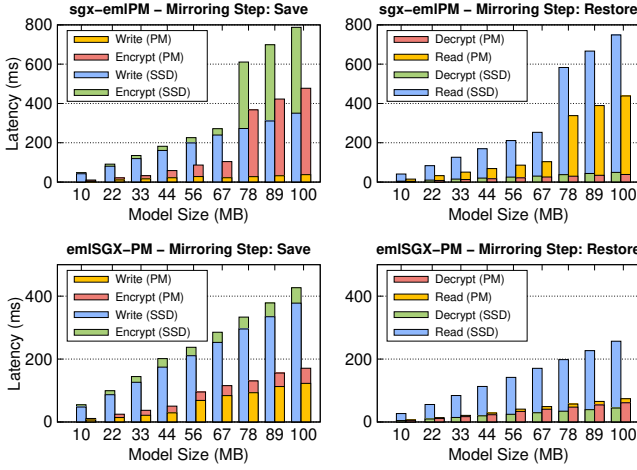


Fig. 7: PM mirroring vs. checkpointing on SSD for *sgx-emIPM* (top) and *emISGX-PM* (bottom).

transactions are $1.6\times$ to $6.9\times$ faster. We justify this behaviour as a result of limited space available for Romulus’ volatile redo log in the SCONE container. These results suggest SGX-ROMULUS is a preferable choice for our ML system, where multiple operations are carried out on persistent models within transactions of relatively larger sizes.

PM mirroring vs. SSD-based checkpointing. Next, we compare the mirroring mechanism in PLINIUS to traditional checkpointing on a SSD using SGX-DARKNET. For SSD checkpointing, we use `ocalls` to `fread` and `fwrite` libC routines to read/write from/to SSD. After each call to `fwrite`, we flush the libC buffers and issue an `fsync`, to ensure data is actually written to secondary storage. We vary model sizes by increasing the total number of convolutional layers. We measure the times to save/mirror-out (encrypt in the enclave and write to PM) and restore/mirror-in (read from PM into enclave and decrypt), and compare these to SSD-based checkpoint saves (encrypt and write to SSD) and SSD-based checkpoint restores (read from SSD into enclave and decrypt), which are the state-of-the-art methods for fault tolerance. All data points are an average of 5 runs.

Figure 7 represents the results obtained on our two servers. As a general observation, in PLINIUS, in-enclave data encryption contributes more to the save-latency (*i.e.*, mirror-out) when compared to writes to PM. For restores in PLINIUS, reads from PM into enclave memory contribute more to the overall latency. When compared to traditional saves and restores on SSD, our mirroring mechanism gives less overhead.

Table 1a shows a performance breakdown of each mirroring steps for saves and restores in PLINIUS, while Table 1b shows the average performance improvements of our mirroring mechanism when compared to SSD-based checkpointing. To reduce the effect of outliers, we evaluate results beneath and beyond the EPC limit separately. The usable EPC size is 93.5 MB, reached for model size 78 MB, due to the presence of other data structures in enclave memory (*e.g.*, temporary buffers used for encryption) as well as enclave code. We observe (Table 1a) that for saves in a real SGX environment, encryption

contributes more (66.4%) to the overall mirroring latency on average for model sizes beneath 78 MB. This jumps to 92.3% once the EPC limit is crossed. This overhead is due to expensive page swapping between the EPC and regular DRAM by the SGX kernel driver. For restores, reads contribute on average 75% and 91.2% for values beneath and beyond the EPC limit respectively. Similarly, we have a high overhead beyond the EPC limit due to the SGX driver’s page swaps. Our results show in-enclave decryption is relatively cheaper.

For the *emISGX-PM* server, without real SGX hardware (hence no expensive page swaps), the main bottleneck is real PM. We observe (Table 1b) that for server *sgx-emIPM*, writes to PM are on average $7.9\times$ and $9.6\times$ faster when compared to writes to SSD for enclave sizes beneath and beyond the EPC limit respectively. SSD writes are generally more expensive due to the expensive `ocalls` and serialization operations to secondary storage. Saves are overall $3.5\times$ and $1.7\times$ faster for enclave sizes beneath and beyond the EPC limit respectively. Similarly, for restores, reads from PM into enclave memory are on average $3\times$ and $1.8\times$ faster for enclave sizes beneath and beyond the EPC limit respectively, when compared to the SSD-based counterpart. Restores are overall $2.5\times$ and $1.7\times$ faster for enclave sizes beneath and beyond the EPC limit. A similar breakdown is done for the *emISGX-PM* node.

Training larger models. Our results suggest PLINIUS is best suited for models with sizes beneath the EPC limit. Models larger than the EPC limit can be trained with PLINIUS but this leads to a significant drop in training performance due to the extensive page swaps by the SGX kernel driver. Figure 7 shows our mirroring mechanism still performs better than SSD-based checkpointing for model sizes beyond the EPC limit. A possible strategy to overcome the EPC limitation could be to distribute the training job over multiple secure CPUs. We will explore this idea in the future. Also, a recent processor release by Intel expands the EPC to 256 MB [4]. This paves the way for applications that leverage PLINIUS to train much larger models more efficiently.

Mirroring frequency. By default PLINIUS does mirroring after every iteration. The mirroring frequency can be easily increased or decreased. All things being equal, a training environment with a small or high frequency of failures will require respectively, small or high mirroring frequencies to achieve good fault tolerance guarantees.

Overhead of data batch decryptions. For efficiency reasons, ML algorithms (*e.g.*, SGD) manipulate training data

(a) Breakdown of mirroring steps (%)			(b) PLINIUS speed-ups		
Save	SGX-emIPM	emISGX-PM	Save	SGX-emIPM	emISGX-PM
Encrypt	66.4%	30.3%	Write	$7.9\times$	$4.5\times$
	92.3%			$9.6\times$	
Write	$33.6\times$	$69.7\times$	Total	$3.5\times$	$3.2\times$
	$7.7\times$			$1.7\times$	
Restore	A	B	Restore	A	B
Read	75%	17.8%	Read	$3\times$	$16.8\times$
	91.2%			$1.8\times$	
Decrypt	25%	82.2%	Total	$2.5\times$	$3.7\times$
	8.8%			$1.7\times$	

TABLE I: Shaded cells: values beyond the EPC size.

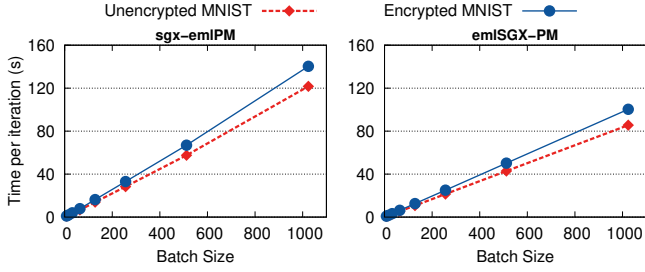


Fig. 8: Variation of iteration times with different batch sizes for encrypted and unencrypted MNIST data.

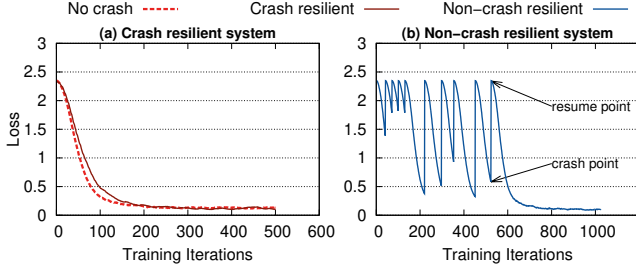


Fig. 9: Crash/resumes are done by randomly killing and restarting the training process every 10 to 15 minutes during model training.

in batches for each training iteration. In this experiment we study the performance impact on total iteration time of batch decryptions of training data into enclave memory. We proceed by comparing the iteration times with different batch sizes for a model being trained via the PLINIUS mechanism, to a model trained with batches of unencrypted data on PM. We recall that in PLINIUS, batches of encrypted training data are read from PM and decrypted in enclave memory for each iteration. All models have 5 LReLU-convolutional layers.

Figure 8 shows the results obtained on both systems. We observe that iterations with batch decryption of data into enclave memory are $1.2\times$ slower on average for both systems. We consider this a relatively small price to pay for data confidentiality during training.

Crash resilience. The main purpose of our experiments here is to demonstrate that PLINIUS’s mirroring mechanism is *crash resilient* (or failure transparent), as well as demonstrate the performance impact on the training process of a *non-crash-resilient* system. We define a crash-resilient system as one capable of recovering its state (*i.e.*, learned parameters) prior to a system crash. The experiments consider models with 5 LReLU-convolutional layers, trained with the MNIST dataset for 500 iterations. We study the variation of the loss while doing random crashes during model training.

Figure 9 presents the results obtained on the *emlSGX-PM* server, but similar results are obtained on *sgx-emIPM*. We proceed by training a model using PLINIUS with 9 random crashes (and resumptions) during the training process. We compare the loss curve obtained here to one obtained without any crashes (baseline). Figure 9(a) shows that despite the crashes, the loss curve follows closely (no breaks at crash and resume points) the one obtained without crashes. This indicates the model parameters are saved and restored correctly using the mirroring mechanism in PLINIUS. In comparison, Figure 9(b)

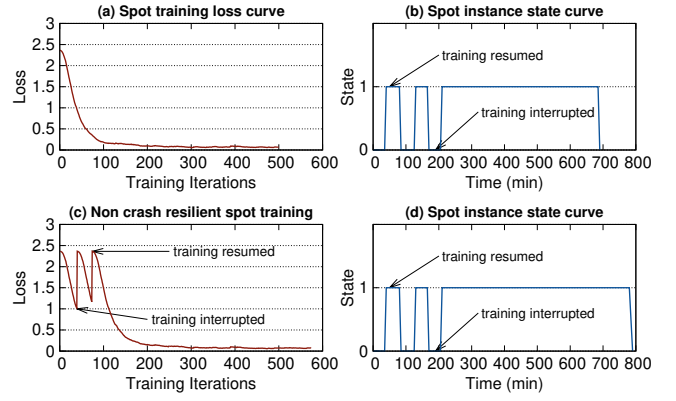


Fig. 10: Model training with AWS EC2 spot instance traces.

shows the loss curve obtained when the system cannot recover its learned parameters following random crashes. For this experiment we run our system while disabling model’s weights saving via our mirroring mechanism. At every resumption point, the model begins the learning process with initial randomized weights, and thus still requires 500 iterations to be fully trained, hence increasing the total iterations (from when training first began) required to train the model to over 1000 in this experiment. This shows the benefit of crash-resilience in an ML system. In the next section, we use a more realistic crash/resume pattern (spot instance trace) to show crash resilience in PLINIUS.

PLINIUS on AWS EC2 Spot instances. A practical use case for PLINIUS framework would be model training on spot instances, such as those offered by Amazon EC2 and Microsoft Azure. Spot instances are liable to many interruptions during their lifetimes, and model training in such a scenario requires efficient fault tolerance guarantees (such as those provided by PLINIUS) to reduce cost and increase efficiency of the training process. We use Amazon EC2 spot instance traces from [38] to simulate a realistic model training scenario with PLINIUS on a spot instance. The spot traces contain market prices of spot instances at different timestamps (5 minutes intervals). To simulate spot model training, we set a *maximum bid price* in our simulator script, and our simulation algorithm periodically (every 5 minutes) compares the *market price* at each timestamp in the spot trace to our bid price. If $max_bid > market_price$, our training process is launched (or continues if it was already running). Otherwise, the training process is killed. We train a model with 12 LReLU-convolutional layers for 500 iterations on server *emlSGX-PM*.

Figure 10(a) shows the loss curve obtained after 500 iterations. As explained in the previous section, this shows PLINIUS is crash resilient as training resumes where it left off prior to the training process being stopped. Figure 10(b) shows a “state curve” of the training process (or spot instance) throughout the training process. The process state is 1 when it is running and 0 otherwise. We observe only 2 interruptions of the training process with our simulation parameters (*i.e.*, maximum bid price of 0.0955). The chosen maximum bid price and spot market price variations will dictate the total number of

interruptions of the spot instance, and hence the total training time (interruption times included). The spot traces used and our simulation scripts are available in the PLINIUS repository.

Figure 10(c) shows us the loss curve obtained when there is no crash resilience (*i.e.*, the model’s state is not saved). With the given simulation parameters (*i.e.*, maximum bid price of 0.0955), there are two interruptions during the training process. As explained in the previous section, it needs to resume training afresh, and hence the combined number of iterations (and total time) from when training first began is increased when compared to its crash resilient counterpart. This further justifies the need for fault tolerance guarantees in such ML scenarios.

CPU and memory overhead. Our mirroring mechanism uses 140 bytes of PM for encryption metadata per layer. The MAC is 16 B, the IV is 12 B, giving 28 B per encrypted parameter buffer. Each layer contains 5 parameter matrices, hence $28 \times 5 = 140$ B per layer. With a model of N layers, we account for $N \times 140$ extra bytes on PM for encryption metadata, small compared to the size of actual models (order of few MBs). The training algorithm is a fairly intensive single-threaded application and it uses 98-100% of the CPU during execution.

Secure inference. PLINIUS can also be used for secure inference. We trained a CNN model with 12 LReLU convolutional layers on the MNIST training dataset, and used the trained model to classify 10’000 grayscale images of handwritten digits in the range $[0 - 9]$. The model (available in the PLINIUS repository) achieved an accuracy of 98.52% with the given hyper-parameters.

GPU and TPU support. Hardware accelerators like Graphics Processing Units (GPUs) and Tensor Processing Units (TPUs) are increasingly used in ML applications. However the former do not integrate TEE capabilities. Recent works like HIX [23], Graviton [37], and Slalom [33] propose techniques to securely offload expensive ML computations to GPUs. Using Darknet’s CUDA extensions, PLINIUS can leverage such techniques to improve training performance. The trained model weights can be securely copied between the secure CPU and the GPU (or TPU) and our mirroring mechanism applied without much changes. We are exploring possible improvements of PLINIUS in this direction.

VII. RELATED WORK

TEE-based schemes. There exists several solutions leveraging trusted hardware (*i.e.*, Intel SGX) for secure ML. Slalom [33] is a framework for secure DNN inference in TEEs. It outsources costly neural network operations to a faster, but untrusted GPU during inference. Occlumency [26] leverages Intel SGX to preserve confidentiality and integrity of user-data during deep learning inference in untrusted cloud infrastructure. Privado [15] implements a secure inference-as-a-service, by eliminating input-dependent access patterns from ML code, hence reducing data leakage risks in the enclave. Chiron [19] leverages Intel SGX for secure ML-as-a-service which prevents disclosure of both data and code.

These systems leverage TEEs for model inference, but without any support for failure recovery. PLINIUS provides a full framework that supports both in-enclave model training and inference with efficient fault tolerance guarantees on PM.

SecureTF [25] integrates TensorFlow ML library for model training and inference in secure SCONE containers. This requires the full TensorFlow library (over 2.5 million LOC [5]) to run inside SGX enclaves, which by design increases the TCB. On the other hand, the trusted portion of PLINIUS comprises 15’900 LOC. The reduction in TCB in PLINIUS when compared to SecureTF is quite obvious; this is better from a security perspective.

Homomorphic encryption (HE)-based schemes. Without trusted hardware enclaves, many privacy-preserving ML methods achieve security via HE-based techniques. HE schemes compute directly over encrypted data. CryptoNets [14] implements inference over encrypted data for pre-trained neural networks. Solutions exist [18] to train and do inference on neural network models using HE.

While these methods ensure privacy of sensitive training and classification data during model training and inference, they have significant performance overhead (up to $1000\times$ slower than TEE-based schemes [20]). PLINIUS provides an orthogonal approach to tackle security, combining Intel SGX enclaves to ensure confidentiality and integrity of models and data sets during training and inference at a much lower cost.

Fault tolerance in ML. A common technique for fault tolerance in ML learning frameworks is checkpointing (restoring) of the model’s state to (from) secondary storage during training (recovery). Several frameworks (*i.e.*, Tensorflow [6], Caffe [24], Darknet [3], *etc.*) rely on secondary storage as persistent storage for training data throughout the training process. Distributing training across several compute nodes improves scalability while increasing fault tolerance.

The above mentioned techniques have huge performance overhead, due to high access times of secondary storage. Following a crash, entire data sets and models must be reloaded into main memory from secondary storage. PLINIUS’s novel mirroring mechanism leverages PM for fault tolerance: upon a crash, the model and the associated training data are readily available in memory. Our design completely obviates the need for expensive serialization (deserialization) of models to (from) secondary storage, and proposes a more efficient approach for handling large amounts of training data.

VIII. CONCLUSION

PLINIUS is the first secure ML framework to leverage Intel SGX for secure model training and PM for fault tolerance. Our novel mirroring mechanism creates encrypted mirror copies of enclave ML models in PM, which are synchronized across training iterations. Our design leverages PM to store byte-addressable training data, completely circumventing expensive disk I/O operations in the event of a system failure. The evaluation of PLINIUS shows that its design substantially reduces the TCB when compared to a system with unmodified libraries, and the mirroring mechanism outperforms disk-based

checkpointing systems while ensuring the system’s robustness upon system failures. Using real-world datasets for image recognition, we show that PLINIUS offers a practical solution to securely train ML models in TEEs integrated with PM hardware at a reasonable cost.

We will extend this work along the following directions. First, we intend to explore GPUs and TPUs by offloading expensive enclave operations on the former without a loss in confidentiality. The extent to which this can be done while preserving confidentiality of the model parameters and training or inference data will be the key area for future work. Second, we wish to explore distributed training using PLINIUS to overcome the SGX EPC limitation. Lastly, we plan to better exploit system parallelism to improve the performance of PLINIUS. This entails redesigning SGX-DARKNET to efficiently support parallel training with threads spawned in the untrusted runtime.

ACKNOWLEDGMENT

This work received funds from the Swiss National Science Foundation (FNS) under project PersIST (no. 178822).

REFERENCES

- [1] “Optimize Your Cloud and Enable Azure Customers to Innovate at Scale with Intel® Optane™ Persistent Memory,” <https://www.intel.la/content/www/xl/es/now/microsoft-azure-optane-innovation-editorial.html>, accessed: Mar 9, 2021.
- [2] “The MNIST Database of Handwritten Digits,” <http://yann.lecun.com/exdb/mnist/>, accessed: May 7, 2020.
- [3] “Darknet: Open Source Neural Networks in C,” <https://pjreddie.com/darknet/>, 2013-2016, accessed: May 7, 2020.
- [4] “Intel SGX Evolves for Data Center,” <https://itpeernetwork.intel.com/intel-sgx-data-center/>, 2019, accessed: Dec 11, 2020.
- [5] “The TensorFlow Open Source Project on Open Hub,” <https://www.openhub.net/p/tensorflow>, 2019, accessed: Dec 11, 2020.
- [6] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard *et al.*, “Tensorflow: A system for large-scale machine learning,” in *USENIX OSDI 2016*.
- [7] M. Al-Rubaie and J. M. Chang, “Privacy-preserving machine learning: Threats and solutions,” *IEEE Security Privacy*, 2019.
- [8] S. Arnavot, B. Trach, F. Gregor, T. Knauth, A. Martin, C. Priebe, J. Lind, D. Muthukumaran, D. O’Keeffe, M. L. Stillwell, D. Goltzsche, D. Eysers, R. Kapitza, P. Pietzuch, and C. Fetzer, “SCONE: Secure linux containers with intel SGX,” in *USENIX OSDI 16*.
- [9] J. V. Bulck, M. Minkin, O. Weisse, D. Genkin, B. Kasikci, F. Piessens, M. Silberstein, T. F. Wenisch, Y. Yarom, and R. Strackx, “Foresadow: Extracting the keys to the intel SGX kingdom with transient out-of-order execution,” in *USENIX Security 2018*.
- [10] I. Corporation, “Intel® Software Guard Extensions Developer Reference for Linux* OS,” https://download.01.org/intel-sgx/sgx-linux/2.8/docs/Intel_SGX_Developer_Reference_Linux_2.8_Open_Source.pdf, 2019.
- [11] A. Correia, P. Felber, and P. Ramalhe, “Romulus: Efficient algorithms for persistent transactional memory,” in *SPAA’18*.
- [12] V. Costan and S. Devadas, “Intel SGX explained,” 2016.
- [13] M. J. Dworkin, “Recommendation for block cipher modes of operation: Galois/counter mode (gcm) and gmac,” Tech. Rep., 2007.
- [14] R. Gilad-Bachrach, N. Dowlin, K. Laine, K. Lauter, M. Naehrig, and J. Wernsing, “Cryptonets: Applying neural networks to encrypted data with high throughput and accuracy,” in *ICML’16*.
- [15] K. Grover, S. Tople, S. Shinde, R. Bhagwan, and R. Ramjee, “Privado: Practical and Secure DNN Inference with Enclaves,” *arXiv preprint arXiv:1810.00602*, 2018.
- [16] D. Gruss, J. Lettner, F. Schuster, O. Ohrimenko, I. Haller, and M. Costa, “Strong and efficient cache side-channel protection using hardware transactional memory,” in *USENIX Security 2017*.
- [17] L. Hanzlik, Y. Zhang, K. Grosse, A. Salem, M. Augustin, M. Backes, and M. Fritz, “Mlcapsule: Guarded offline deployment of machine learning as a service,” *arXiv preprint arXiv:1808.00590*, 2018.
- [18] E. Hesamifard, H. Takabi, M. Ghasemi, and C. Jones, “Privacy-preserving machine learning in cloud,” in *Proceedings of the 2017 on Cloud Computing Security Workshop*, 2017.
- [19] T. Hunt, C. Song, R. Shokri, V. Shmatikov, and E. Witchel, “Chiron: Privacy-preserving machine learning as a service,” *arXiv preprint arXiv:1803.05961*, 2018.
- [20] N. Hynes, R. Cheng, and D. Song, “Efficient deep learning on multi-source private data,” *arXiv preprint arXiv:1807.06689*, 2018.
- [21] C. Iorgulescu, R. Azimi, Y. Kwon, S. Elnikety, M. Syamala, V. R. Narasayya, H. Herodotou, P. Tomita, A. Chen, J. Zhang, and J. Wang, “Perfiso: Performance isolation for commercial latency-sensitive services,” in *USENIX ATC 2018*.
- [22] J. Izraelevitz, J. Yang, L. Zhang, J. Kim, X. Liu, A. Memaripour, Y. J. Soh, Z. Wang, Y. Xu, S. R. Dulloor *et al.*, “Basic performance measurements of the intel optane dc persistent memory module,” *arXiv preprint arXiv:1903.05714*, 2019.
- [23] I. Jang, A. Tang, T. Kim, S. Sethumadhavan, and J. Huh, “Heterogeneous isolated execution for commodity gpus,” in *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*, ser. ASPLOS ’19. New York, NY, USA: Association for Computing Machinery, 2019, p. 455–468. [Online]. Available: <https://doi.org/10.1145/3297858.3304021>
- [24] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, “Caffe: Convolutional architecture for fast feature embedding,” in *ACM Multimedia 2014*.
- [25] R. Kunkel, D. L. Quoc, F. Gregor, S. Arnavot, P. Bhatotia, and C. Fetzer, “SecureTF: A Secure TensorFlow Framework,” in *Middleware’20*.
- [26] T. Lee, Z. Lin, S. Pushp, C. Li, Y. Liu, Y. Lee, F. Xu, C. Xu, L. Zhang, and J. Song, “Occlumency: Privacy-preserving remote deep-learning inference using SGX,” in *MobiCom’19*.
- [27] L. Lersch, X. Hao, I. Oukid, T. Wang, and T. Willhalm, “Evaluating persistent memory range indexes,” *VLDB Endowment 2019*.
- [28] E. Liberty, Z. Karnin, B. Xiang, L. Rouesnel, B. Coskun, R. Nallapati, J. Delgado, A. Sadoughi, Y. Astashonok, P. Das *et al.*, “Elastic Machine Learning Algorithms in Amazon SageMaker,” *SIGMOD’20*.
- [29] P. Mohassel and Y. Zhang, “SecureML: A system for scalable privacy-preserving machine learning,” in *IEEE S&P 2017*.
- [30] O. Oleksenko, B. Trach, R. Krahn, M. Silberstein, and C. Fetzer, “Varys: Protecting SGX enclaves from practical side-channel attacks,” in *USENIX ATC 2018*.
- [31] M. Schwarz, S. Weiser, D. Gruss, C. Maurice, and S. Mangard, “Malware guard extension: Using SGX to conceal cache attacks,” in *International Conference on Detection of Intrusions and Malware, and Vulnerability Assessment*, 2017.
- [32] Y. Shen, H. Tian, Y. Chen, K. Chen, R. Wang, Y. Xu, Y. Xia, and S. Yan, “Occlum: Secure and Efficient Multitasking Inside a Single Enclave of Intel SGX,” ser. ASPLOS’20.
- [33] F. Tramer and D. Boneh, “Slalom: Fast, verifiable and private execution of neural networks in trusted hardware,” in *ICLR’19*.
- [34] F. Tramèr, F. Zhang, A. Juels, M. K. Reiter, and T. Ristenpart, “Stealing machine learning models via prediction apis,” in *USENIX Security 2016*.
- [35] C. Tsai, D. E. Porter, and M. Vij, “Graphene-sgx: A practical library OS for unmodified applications on SGX,” in *USENIX ATC 2017*.
- [36] H. Volos, A. J. Tack, and M. M. Swift, “Mnemosyne: Lightweight persistent memory,” in *ASPLOS’11*.
- [37] S. Volos, K. Vaswani, and R. Bruno, “Graviton: Trusted execution environments on gpus,” in *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*. Carlsbad, CA: USENIX Association, Oct. 2018, pp. 681–696. [Online]. Available: <https://www.usenix.org/conference/osdi18/presentation/volos>
- [38] C. Wang, Q. Liang, and B. Urgaonkar, “An empirical analysis of amazon EC2 Spot Instance features affecting cost-effective resource procurement,” *TOMPECS’18*.
- [39] N. Weichbrodt, P. Aublin, and R. Kapitza, “sgx-perf: A performance analysis tool for intel SGX enclaves,” in *Middleware’18*.
- [40] J. Yang, J. Kim, M. Hoseinzadeh, J. Izraelevitz, and S. Swanson, “An empirical guide to the behavior and use of scalable persistent memory,” in *USENIX FAST 2020*.
- [41] P. Zuo and Y. Hua, “Secpm: a secure and persistent memory system for non-volatile memory,” in *HotStorage’18*.