

# Understanding and Defending against the Security Threats on Mobile and IoT Devices with Hardware-backed Trusted Computing Primitives

PHD THESIS PROPOSAL

ZHICHUANG SUN

NORTHEASTERN UNIVERSITY  
KHOURY COLLEGE OF COMPUTER SCIENCES

PHD COMMITTEE

Long Lu, Northeastern University  
Engin Kirda, Northeastern University  
Guevara Noubir, Northeastern University  
Somesh Jha, University of Wisconsin-Madison

January 18, 2021

## Abstract

Mobile devices have become an essential part of our daily life while IoT devices are being rapidly deployed in our home, factories, and infrastructures. These omnipresent devices bring in great convenience along with growing security and privacy concerns. For example, the attacks on IoT devices are becoming increasingly complex and destructive. New attack surfaces are also being opened up on mobile devices along with the wide deployment of machine learning (ML) technology, *e.g.*, the privacy of on-device machine learning models.

In this thesis, I present two novel systems and one large-scale study to understand and combat the newly emerging attacks on both IoT devices and mobile devices. The security of these two systems is built on hardware-backed trusted computing primitives to provide strong security guarantee. First, I present OAT, an attestation system that captures the integrity of both control flow and critical data involved in an operation execution to detect advanced attacks on IoT devices. OAT advances the state-of-the-art attestation on IoT devices by providing strong control-flow verification and light-weight critical data integrity check. Second, I present the first large-scale study of insufficient model protection in mobile apps to uncover this less understood problem of model privacy on mobile devices. Based on the 46,753 trending Android apps collected from the US and Chinese app markets, alarmingly, this study shows that 41% of ML apps do not protect their models at all. For those apps that use model protection or encryption, we were able to extract the models from 66% of them via unsophisticated dynamic analysis techniques. Our financial and security impact analysis on (stolen) models indicates that the consequence can be severe for both the model vendors and the app owners. Third, to protect on-device ML models, I propose ShadowNet, a secure and efficient model inference system built to protect model privacy on mobile devices. ShadowNet offers a novel idea to outsource the heavy part of the model inference to the untrusted world (including GPU) for acceleration without leaking the model weights. The security of ShadowNet is rooted in the TEE, which can protect the model privacy even when the OS is compromised.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Problem Statement . . . . .	3
1.2	Thesis Statement . . . . .	3
1.3	Approaches Overview . . . . .	4
<b>2</b>	<b>Background and Related Work</b>	<b>4</b>
2.1	Attacks on IoT Devices and Backends . . . . .	4
2.2	Arm TrustZone . . . . .	5
2.3	Remote Attestation . . . . .	5
2.4	On-device Machine Learning . . . . .	5
2.5	Secure Machine Learning . . . . .	6
<b>3</b>	<b>Attesting Operation Integrity of IoT Devices</b>	<b>6</b>
3.1	Introduction of the Attacks on IoT Devices . . . . .	6
3.2	Methodology . . . . .	7
3.3	Evaluation Results . . . . .	8
<b>4</b>	<b>A Large-scale Study on Insufficient Machine Learning Model Protection in Mobile Apps</b>	<b>9</b>
4.1	Analysis Overview . . . . .	10
4.2	Q1: How Widely Is Model Protection Used in Apps? . . . . .	10
4.2.1	Methodology . . . . .	11
4.2.2	Findings and Insights . . . . .	11
4.3	Q2: How Robust Are Existing Model Protection Techniques? . . . . .	13
4.3.1	Methodology . . . . .	13
4.3.2	Findings and Insights . . . . .	14
4.4	Q3: What Impacts can (Stolen) Models Incur? . . . . .	15
4.4.1	Financial Impact . . . . .	15
4.4.2	Security Impact . . . . .	15
<b>5</b>	<b>Research Plan</b>	<b>16</b>
5.1	ShadowNet: a Secure and Efficient Model Inference Scheme . . . . .	16
5.2	Formal Security Analysis of the Scheme . . . . .	16
5.3	System Design Challenges . . . . .	17
5.4	Milestones . . . . .	17

# 1 Introduction

## 1.1 Problem Statement

Mobile and IoT devices are omnipresent. Our shopping, social networking, entertainment, work, and payment are made much more convenient with mobile phones. IoT devices are also being rapidly deployed in home automation, industrial automation, smart cities, agriculture and so on, making our home more smart and our industry more efficient. However, with great convenience also comes great concerns about the security and privacy of our digital assets. First, the attacks on IoT devices are evolving and become more and more complex and destructive. Second, with the newly emerging technologies like AI being deployed on mobile devices, new attack surface are being opened up for the attackers, *e.g.*, the privacy of machine learning models on our mobile phones.

**Advanced attacks on IoT devices:** Internet-of-Things (IoT) and Cyber-Physical Systems (CPS) are being rapidly deployed in smart homes, automated factories, intelligent cities, and more. As a result, embedded devices, playing the central roles as sensors, actuators, or edge-computing nodes in IoT systems, are becoming attractive targets for cyber attacks. Unlike computers, attacks on embedded devices can cause not only software failures or data breaches but also physical damage. Moreover, a compromised device can trick or manipulate the IoT backend (*e.g.*, remote controllers or in-cloud services): hijacking operations and forging data.

The early attacks on IoT devices usually abuse the misconfiguration of IoT devices, like default passwords [73], which can be mitigated through education on the users and security-aware default configuration. Nowadays, more powerful attacks that exploit the implementation bug of wireless protocol have been demonstrated practical and destructive [61]. Return-oriented programming (ROP) and data-only attacks are also proven to be easy to launch on embedded devices, as demonstrated on vulnerable industrial robot controllers [59].

Unfortunately, today’s IoT backends cannot protect themselves from manipulations by compromised IoT devices. This is due to the lack of a technique for remotely verifying if an operation performed by an IoT device has been disrupted, or any critical data has been corrupted while being processed on the device. As a result, IoT backends are forced to blindly trust remote devices for faithfully performing assigned operations and providing genuine data. Our work aims to make this trust verifiable, and therefore, prevent compromised IoT devices from deceiving or manipulating the IoT backend.

**Model privacy on mobile devices:** With the rapid development of machine learning technology, mobile app developers have been quickly adopting *on-device machine learning* (ML) techniques to provide artificial intelligence (AI) features, such as facial recognition, augmented/virtual reality, image processing, voice assistant, etc. This trend is now boosted by new AI chips available in the latest smartphones [1], such as Apple’s Bionic neural engine, Huawei’s neural processing unit, and Qualcomm’s AI-optimized SoCs.

Compared to performing ML tasks in the cloud, on-device ML (mostly model inference) offers unique benefits desirable for mobile users as well as app developers. For example, it avoids sending (private) user data to the cloud and does not require network connection. For app developers or ML solution providers, on-device ML greatly reduces the computation load on their servers.

On-device ML inference inevitably stores ML models locally on user devices, which however creates a new security challenge. Commercial ML models used in apps are often part of the core intellectual property (IP) of vendors. Such models may fall victim to theft or abuse, if not sufficiently protected. In fact, on-device ML makes model protection much more challenging than server-side ML because models are now stored on user devices, which are fundamentally untrustworthy and may leak models to curious or malicious parties.

The consequences of model leakage are quite severe. First, with a leaked model goes away the R&D investment of the model owner, which often includes human, data, and computing costs. Second, when a proprietary model is obtained by unethical competitors, the model owner loses the competitive edge or pricing advantage for its products. Third, a leaked model facilitates malicious actors to find adversarial inputs to bypass or confuse the ML systems, which can lead to not only reputation damages to the vendor but also critical failures in their products (*e.g.*, fingerprint recognition bypass).

## 1.2 Thesis Statement

Driven by the great security and privacy concern about the omnipresent mobile and IoT devices, on which our life has growing dependence, my research aims to : (1) secure the IoT devices in the face of evolving attacking techniques; (2) advance the understanding of the newly emerging attack surface on mobile devices, namely the privacy problem of on-device ML models, and (3) design novel systems to protect the model privacy on mobile devices.

For IoT devices, we need to improve our attack detection capabilities. In addition to detecting traditional attacks that change the firmware, we should also develop new techniques to capture advanced attacks that can manipulate the device behavior without modifying the firmware, *e.g.*, control-flow hijacking and data-only attacks.

For mobile devices, we need to systematically study the new problem of ML model privacy, figure out how widely the problem is and what impacts it can incur, and understand the challenges of protecting on-device ML models.

Once we get a understanding of the new security threats on mobile devices, we should also attempt to design a secure and efficient way to run model inference on mobile devices so as to protect the ML models' privacy without breaking existing ML applications.

### 1.3 Approaches Overview

In this thesis, I advance the state-of-the-art research on the security of IoT devices and mobile devices with novel security system based on hardware-backed trusted computing primitives and systematic study of the ML model privacy problem on mobile devices.

First, I present OAT [66], which is the first attestation method that captures both control-flow and data-only attacks on embedded devices. The root of security is built on top of the TEE. Using this method, IoT backends can now verify if a remote device is trustworthy when it claims it has performed an operation, sent in a service request, or transported back data from the field. In addition, unlike traditional attestation methods, which only output a binary result, our method allows verifiers to reconstruct attack execution traces for postmortem analysis.

Second, I present the first large-scale study of ML model protection and theft on mobile devices based on 46,753 trending Android apps collected from the US and the Chinese app markets [67]. Our study aims to shed light on the less understood risks and costs of model leakage/theft in the context of on-device ML. We present our study that answers the following questions with ample empirical evidence and observations: (1) How widely is model protection used in apps? (2) How robust are existing model protection techniques? and (3) What impacts can (stolen) models incur?.

Finally, I propose ShadowNet as part of my thesis, which is a secure and efficient model inference system designed for mobile devices. The security of ShadowNet is also rooted in the TEE. However, ShadowNet avoids the naive approach of moving the whole model inference into the TEE which can easily exhaust the limited resource of the TEE and lose access to the untrusted hardware accelerators. Instead, ShadowNet offers a novel idea based on linear transformation to outsource the heavy linear layers of the model to the untrusted world (including GPU) for acceleration without leaking the model weights. With this novel design, ShadowNet has a potential to achieves hardware-backed security with little overhead and small TCB size.

## 2 Background and Related Work

### 2.1 Attacks on IoT Devices and Backends

Embedded devices, essential for IoT, have been increasingly targeted by powerful attacks. For instance, hackers have managed to subvert different kinds of smart home gadgets, including connected lights [61], locks [42], etc. In industrial systems, robot controllers [59] and PLCs (Programmable Logic Controller) [30] were exploited to perform unintended or harmful operations. The same goes for connected cars [52, 44], drones [39], and medical devices [31, 60]. In addition, large-scale IoT deployments were compromised to form botnets via password cracking [73], and recently, vulnerability exploits [50].

Meanwhile, advanced attacks quickly emerged. Return-oriented Programming (ROP) was demonstrated to be realistic on RISC [22], and particularly Arm [49], which is the common architecture for today's embedded devices. Data-only attacks [27, 43] are not just applicable but well-suited for embedded devices [74], due to the data-intensive or data-driven nature of IoT.

Due to the poor security of today's embedded devices, IoT backends (*e.g.*, remote IoT controllers and in-cloud services) are recommended to operate under the assumption that IoT devices in the field can be compromised and should not be fully trusted [65]. However, in reality, IoT backends are often helpless when deciding whether or to what extent it should trust an IoT device. They may resort to the existing remote attestation techniques, but these techniques are only effective at detecting the basic attacks (*e.g.*, device or code modification) while leaving advanced attacks undetected (*e.g.*, ROP, data-only attacks, etc.). As a result, IoT backends have no choice but to trust IoT devices and assume they would faithfully execute commands and generate genuine data or requests. This blind and unwarranted trust can subject IoT backends to deceptions and manipulations. For example, a compromised

robotic arm can drop a command yet still report a success back to its controller; a compromised industrial syringe can perform an unauthorized chemical injection, or change an authorized injection volume, without the controller’s knowledge.

## 2.2 Arm TrustZone

TrustZone is a hardware feature available on both Cortex-A processors (for mobile and high-end IoT devices) and Cortex-M processors (for low-cost embedded systems). TrustZone renders a so-called “Secure World”, an isolated environment with tagged caches, banked registers, and private memory for securely executing a stack of trusted software, including a tiny OS and trusted applications (TA). In parallel runs the so-called “Normal World”, which contains the regular/untrusted software stack. Code in the Normal World, called client applications (CA), can invoke TAs in the Secure World. A typical use of TrustZone involves a CA requesting a sensitive service from a corresponding TA, such as signing or securely storing a piece of data. In addition to executing critical code, TrustZone can also be used for mediating peripheral access from the Normal World.

## 2.3 Remote Attestation

Early works on remote attestation, such as [70][55], were focused on static code integrity, checking if code running on remote devices has been modified. A series of works [16, 58, 29, 48] studied the Root of Trust for remote attestation, relying on either software-based TCB or hardware-based TPM or PUF. Armknecht et al. [17] built a security framework for software attestation.

Other works went beyond static property attestation. Haldar et al. [71] proposed the verification of some high-level semantic properties for Java programs via an instrumented Java virtual machine. ReDAS [47] verified the dynamic system properties. Compared with our work, these previous systems were not designed to verify control-flow or dynamic data integrity. Further, their designs do not consider bare-metal embedded devices or IoT devices. Some recent remote attestation systems addressed other challenges. A tool called DARPA [45] is resilient to physical attacks. SEDA [19] proposed a swarm attestation scheme scalable to a large group of devices. In contrast, we propose a new remote attestation scheme to solve a different and open problem: IoT backend’s inability to verify if IoT devices faithfully perform operations without being manipulated by advanced attacks (i.e., control-flow hijacks or data-only attacks). Our attestation centers around OEI, a new security property we formulated for bare-metal embedded devices. OEI is operation-oriented and entails both control-flow and critical data integrity.

A recent work called C-FLAT [14] is closely related to our work. It enabled control-flow attestation for embedded devices. However, it suffers from unverifiable hashes, especially when attested programs have nested loops and branches. This is because verifying a control-flow hash produced by C-FLAT requires the knowledge of all legitimate control-flow hashes, which are impossible to completely pre-compute due to the unbounded number of code paths in regular programs (i.e., the path explosion problem). In comparison, OAT uses a new hybrid scheme for attesting control-flows, which allows deterministic and fast verification. Moreover, OAT verifies Critical Variable Integrity and can detect data-only attacks, which C-FLAT and other previous works cannot.

## 2.4 On-device Machine Learning

**The Trend of On-device Machine Learning:** Currently, there are two ways for mobile apps to use ML: cloud-based and on-device. In cloud-based ML, apps send requests to a cloud server, where the ML inference is performed, and then retrieve the results. The drawbacks include requiring constant network connections, unsuitable for real-time ML tasks (e.g., live object detection), and needing raw user data uploaded to the server. Recently, on-device ML inference is quickly gaining popularity thanks to the availability of hardware accelerators on mobile devices and the the ML frameworks optimized for mobile apps. On-device ML avoids the aforementioned drawbacks of cloud-based ML. It works without network connections, performs well in real-time tasks, and seldom needs to send (private) user data off the device. However, with ML inference tasks and ML models moved from cloud to user devices, on-device ML raises a new security challenge to model owners and ML service providers: how to protect the valuable and proprietary ML models now stored and used on user devices that cannot be trusted.

**The Delivery and Protection of On-device Models :** Typically, on-device ML models are trained by app developers or ML service providers on servers with rich computing resources (e.g., GPU clusters and large storage servers). Trained models are shipped with app installation packages. A model can also be downloaded separately after app installation to reduce the app package size. Model inference is performed by apps on user devices, which

relies on model files and ML frameworks (or SDKs). To protect on-device models, some developers encrypt/obfuscate them, or compile them into app code and ship them as stripped binaries[5, 10]. However, such techniques only make it difficult to reverse a model, rather than strictly preventing a model from being stolen or reused.

## 2.5 Secure Machine Learning

Existing research on secure machine learning covers both the end devices and the cloud. Offline Model Guard (OMG) [20] provides a secure model inference framework for mobile device based on SANCTUARY[21], a user space enclave based on Arm TrustZone. OMG allows the model inference framework runs fully inside SANCTUARY enclave to protect model privacy. MLCapsule [41] also deploy the model on the client side to protect the user input from being sent to the untrusted cloud end. At the same time, it runs the model inference inside SGX to prevent the model from being leaked to the client. GPU is not guarded by SANCTUARY and SGX, so both OMG and MLCapsule do not support secure GPU acceleration. DarknetTZ[57] is a secure machine learning framework built on top of Arm TrustZone. It allows a few selected layers to be running inside TEE to protect part of the model. By running heavy linear layers inside TEE, it also poses resource challenges(like memory)on TEE. Comparing with OMG, MLCapsule and DarknetTZ, ShadowNet has a small TCB inside TEE and allows secure outsource of linear layers onto GPU. Graviton [72] proposes TEE extension for GPU hardware, thus allowing GPU tasks like machine learning to be running securely on GPU. It is a promising feature but requires hardware changes on GPU. Secloak [54] partitions GPU into secure world with Arm TrustZone to run GPU tasks securely at high performance penalty.

Research on securing machine learning on the cloud end is an active area. TensorSCONE[51] propose a secure machine learning framework running in the untrusted cloud. TensorSCONE integrates TensorFlow with the secure Linux container technology SCONE[18] guarded by SGX. TF Trusted [12] leverages custom operations to send gRPC messages into the Intel SGX device via Google Asylo[4] where the model is then run by Tensorflow Lite. Running model inference inside TEE faces performance challenges due to limited memory and lack of GPU acceleration. Occlumency [53] provides a suite of heuristic techniques based on Caffe and improves inference speed by 3.6 times. Despite promising, these works do not support GPU acceleration.

YerbaBuena [40] partitions the model into frontnets(like the first layer) and backnets,and execute the frontnets inside SGX to protect the user input from being leaked to the untrusted cloud while running backnets unprotected to leverage hardware acceleration. Slalom [68] splits DNN into linear and nonlinear layers, and outsources linear layers to GPU for acceleration with masked input. It verifies the linear layer’s results and computes the nonlinear layers inside SGX. Slalom protects the user input privacy but not the model weights. SecureNets[28] transforms both input and linear layer’s weights into matrix, and applies matrix transformation proposed in [62] to hide the non-zero elements, then sends them to the untrusted cloud for acceleration. It is not clear whether SecureNets supports depthwise convolution and convolution with stride. ShadowNet does not require transforming input and weights into matrix, and is compatible with existing linear operations.

There are also secure ML with cryptographic approach. CryptoNets[34] applies Homomorphic Encryption(HE) on neural networks and runs model inference on encrypted data to protect user input privacy. Jiang et. al [46] presents a solution to encrypt a matrix homomorphically and perform arithmetic operations on encrypted matrices. It protects both user data and model. TF Encrypted [11] enables training and prediction over encrypted data via secure multi-party computation and homomorphic encryption. SafetyNets[32] designs a Interactive Protocol(IP) that allows clients to verify the correctness of a class of DNNs running on the untrusted cloud by asking for a short mathematical proof.

To ensure the integrity of model weights, Uchida et. al [69] and Zhang et. al [75] embed watermarks into deep neural model parameters, while training the models. DeepAttest [26] encodes fingerprint in DNN weights to prevent weight modification. These works are incapable of preventing weight leakage.

## 3 Attesting Operation Integrity of IoT Devices

### 3.1 Introduction of the Attacks on IoT Devices

IoT devices are being widely deployed in our home, industries and city infrastructures. These omnipresent devices are becoming increasingly connected and mission-critical, allowing the operators to monitor and control these physically scattered devices, *e.g.*, home appliances, industry robots, or air quality sensors, from anywhere, at anytime.

Unfortunately, several characteristics of the IoT devices make them attractive targets for the attackers. First, the universally identical software system. IoT devices are usually designed for certain purpose that can be widely deployed and universally managed without the involvement of individual users. For one type of device, the software systems of different devices are usually identical. As a result, the attack on one device can be easily repeated on another. Second, wide connectivity. To allow universal management from the vendors and easy control from the users, IoT devices are usually connected to a cloud hub, some of them also allows access and control from the end user’s mobile device. Besides, some devices like home appliances, also allow easy access and configuration from other devices under the same local area network, *e.g.*, WiFi. Wide connectivity opens a large attack space for the attackers and give them more power once they take control of an insider node of these connected devices. Third, easy availability. IoT devices are usually cheap and can be bought online. The attackers can simply buy one and reverse engineer the device for fun at low cost[33].

The attacks on IoT devices are also evolving. Earlier large scale attacks usually exploit insecure default settings (*e.g.*, passwords) on IoT devices[73]. Recently, researchers [61] have demonstrated that more advanced attacks on IoT devices, *e.g.*, exploiting the implementation bug of wireless protocol, can be destructive by spreading IoT worms. An experimental security analysis of industrial robot controller [59] also shows that more advanced attacks (*e.g.*, control-flow hijacks and data-only attacks) are not only possible but also easier to launch on those industrial robot controllers due to dated software and libraries.

Previous research on IoT devices focuses on the static property, *e.g.*, code integrity [63, 70, 71, 64, 24, 47]. Code integrity can be used to detect permanent attacks that modifies the firmware, but it can not detect advanced attacks, *e.g.*, control-flow hijacks and data-only attacks, which do not modify the firmware. A recent work C-FLAT [14] took one step forward to introduce control-flow integrity (CFI) attestation. However, the proposed control-flow verification is non-deterministic due to the path explosion issue, and it does not cover data-only attacks. There lacks a general approach for the IoT backends (*e.g.*, IoT hub in the cloud) to detect if any such advanced attacks have happened while an operation is performed on the IoT end devices. As a result, the IoT backends have to blindly trust that the IoT devices perform each operation faithfully without being manipulated.

### 3.2 Methodology

To address this problem, we first formulate a new security property for embedded devices, called “*Operation Execution Integrity*” or *OEI*. We then design and build OAT, a system that enables remote OEI attestation for Arm-based bare-metal embedded devices. OAT can detect both unexpected control-flow and data manipulations in an efficient way, which existing attestation methods cannot check.

Our formulation of OEI is inspired by the observation that the tasks performed on embedded devices are usually composed of several operations. Operation here means a logically independent task performed by an embedded device. For example, taking a measurement of the temperature, moving the robotic arm from position A to position B are two different operations. These operations are critical to the fulfillment of the tasks received from the IoT hub. For simplicity of discussion, we assume an operation always has a single pair of entry and exit.

We use operation as the granularity of attestation to avoid the high overhead of always-on measurements. The operations are triggered only when the device receives new commands from the IoT hub. In this way, OAT allows the attestation to be triggered on need, minimizing the energy consumption or overhead on the end devices when no operation is triggered.

OEI captures the control-flow and critical data integrity within an operation. Satisfying OEI entails an operation is executed without any manipulation of control-flow and critical data. For control-flow verification, we design an efficient measurement scheme which is light-weight and deterministic. It is based on the observation that there are two types of control-flow events: forward and backward. Forward control-flow events refer to branches and function calls, while backward control-flow events refer to function returns. For a normal execution without control-flow hijacks, once the forward control-flow is determined, there is only one legal backward control-flow. OAT uses two different measurements for forward and backward control-flow events. For forward events, we record the taken/not-taken information for conditional jumps, which takes only one bit space to store. We also record the destination addresses of indirect jumps/function calls. We do not record any information for unconditional branches or direct function calls. For backward events, we do not record any destination addresses of function returns. Instead, we calculate a hash over all function return events, which comprise of the start addresses and destination addresses. During verification, our verification engine use abstract execution on the binary with the recorded forward events information to reconstruct the control flow and calculate the hash. We then compare the inferred hash with the hash measured on the device. A hash mismatch means control-flow hijacks happened during measurement. The use of hash for backward control-flow events greatly reduced the size of the control-flow measurements.

Table 1: Runtime overhead measured on 5 real embedded programs

Prog.	Operation Exec. Time			OAT Instrumentation Statistics					Blob Size (B)	Verification Time (s)
	w/o OEI (s)	w/ OEI (s)	Overhead (%)	B.Cond	Def-Use	Ret	Icall/Ijmp	Critical Var.		
SP	10.19	10.38	1.9%	488	2	1946	1	20	69	5.6
HA	5.28	5.36	1.6%	147	91	33	2	6	44	0.61
RM	10.01	10.13	1.3%	901	100	100	100	7	913	1.74
RC	2.55	2.66	4.5%	14	33	1	1	8	10	0.13
LC	5.33	5.56	4.4%	931	2420	10	10	4	205	1.35
Avg.	N/A	N/A	2.7%	496	529	418	23	9	248	1.89

Table 2: Number of Instrumentation Sites: Value-based (R1) and Address-based (R2)

	SP	HA	RM	RC	LC	Avg.
R1	56	37	57	20	41	-
R2	140	388	842	45	131	-
R1 / R2	40%	9.5%	6.8%	44.4%	31.2%	26%

Our critical data integrity is also customized for embedded device to reduce the overhead. By separating critical data from non-critical data, and only enforcing integrity check for critical data, we can reduce the code size and runtime overhead. Critical data refers to those critical variables that affects the execution of an operation. We define two types of critical variables: semantically critical variables and conditional variables. Semantically critical variables need to be annotated by the programmer. For example, for operations involve sensors, the variables that store the value read from sensors are semantically critical. For operations involve the robotic arm movement, the variables that store the actual distance or angle of the arm movement are semantically critical. Conditional variables refers to variables used in the condition expression of a control-flow transfer sentence, *e.g.*, variables used in `if` or `while` condition expressions. Conditional variables can be automatically annotated by compiler. Manipulating these variables will bend the control-flow to favor the attackers.

Our critical variable integrity check is different from previous works on data integrity check [25, 15, 23] as they require heavy instrumentation. Previous works like DFI [25] and DataShield [23] performs address-based check. They need to instrument every write instruction to enforce destination address check, no matter it is protecting all program data or selected part of program data. Our goal is to enforce the integrity check of critical variables. We design a value-based integrity check, which only needs to instrument those instructions that need to access those critical variables. Our critical variable integrity requires the read value of each critical variable matches the write value from the last legal write instruction. We put more effort on compiler side to statically identify all legal read/write access instructions to critical variables, and instrument these relevant instructions. When the ratio of critical variables and all variables is low, our approach incurs much less instrumentation overhead comparing with address-based checks.

### 3.3 Evaluation Results

We evaluate the end-to-end overhead of OAT on five open-source embedded programs. These five programs represent a reasonable level of variety. The 5 selected embedded programs are:

- *Syringe Pump* (SP) is a remotely controlled liquid-injection device, often used in healthcare and food processing settings. We apply OEI attestation to the “injection-upon-command” operation.
- *House Alarm System* (HA) [35] is an IoT device that, when user-specified conditions are met, takes a picture and triggers an alarm. We apply OEI attestation to its “check-then-alarm” operation.
- *Remote Movement Controller* (RM) [37] is an embedded device that allows the physical movement of its host to be controlled remotely. We attest the “receive-execute-command” operation.
- *Rover Controller* (RC) [38] controls the motor on a rover. We attest the “receive-execute-command” operation.
- *Light Controller* (LC) [36] is a smart lighting controller. We attest the “turn-on/off-upon-command” operation.

**Compile-time Overhead:** Embedded programs are usually sensitive to program size. We measured the code size and compilation time increase for the five test programs. Figure 1 shows the results for each program. The

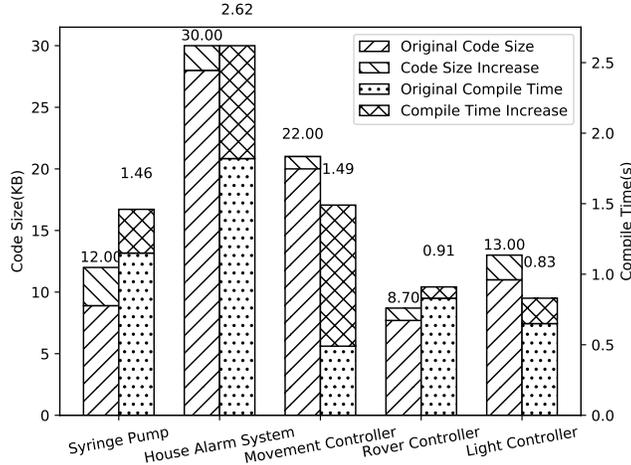


Figure 1: Compile-time overhead

Table 3: Control-flow Trace Size (Bytes): With Return Hash (R1) and Without Return Hash (R2)

	SP	HA	RM	RC	LC	Avg.
R1	69	44	913	10	205	-
R2	42941	3772	13713	585	13725	-
R1 / R2	0.2%	1.1%	6.7%	1.7%	1.5%	2.24%

absolute code size increase is less than 3 KB, which is acceptable for embedded devices. The compilation delay caused by extra code analysis and instrumentation is less than 1 second for the tested programs. Considering that the embedded programs are usually small by size, the compilation delay will not significantly affect the developer’s compilation experience.

**Operation Execution Time & Instrumentation Statistics:** For each test program, we also measure the runtime overhead caused by attestation. We measure the operation execution time with and without OEI attestation. The results are shown in the column of “Operation Exec. Time” in Table 1. The sub-column “Overhead” shows the relative delay caused by OAT to operation executions, averaging 2.7%. As we can see, the execution delay caused by our attestation may vary due to the length/duration of the attested operation and the frequency of critical variable def-use events. We also show the statistics of control-flow and critical variables in Table 1.

**Value-based Check vs. Addressed-based Check:** To compare the difference between value-based check and address-based check, we measured the number of instrumented instructions needed in both cases for all of the test programs. As shown in Table 2, on average, CVI’s instrumentation is 74% less than the instrumentation required by address-based checking (i.e., a 74% reduction).

**Space-efficiency of Hybrid Attestation:** Our control-flow attestation uses the hybrid scheme consisting both forward traces and backward hashes. We compared the sizes of the control-flow traces produced by OAT (R1 in Table 3) and the traces produced by pure trace-based CFI (R2 in Table 3). On average, OAT’s traces take only 2.24% of space as needed by control-flow traces (i.e., a 97% reduction). It shows that our hybrid scheme is much better suited for embedded devices in terms of space-efficiency.

## 4 A Large-scale Study on Insufficient Machine Learning Model Protection in Mobile Apps

Machine Learning (ML) technology is being rapidly developed and deployed to bring intelligence to our life. Specifically, the breakthrough of Deep Learning (DL) in image processing and NLP has unlocked many popular product features like face recognition and chatting robots. In the beginning, popular ML services provided to the end users rely on heavy work loads running in the cloud due to its high demand for computing resources. Nowadays,

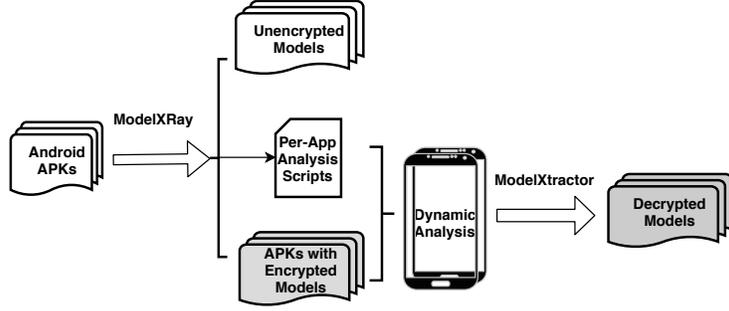


Figure 2: Overview of Static-Dynamic App Analysis Pipeline

with new AI chips being available in the latest smartphones [1], such as Apple’s Bionic neural engine, Huawei’s neural processing unit, and Qualcomm’s AI-optimized SoCs, on-device machine learning (ML) is quickly gaining popularity among mobile apps. It greatly saved the latency of sending data between end devices and the cloud, and allows offline model inference while preserving user privacy. However, ML models, considered as core intellectual properties of model owners, are now stored on billions of untrusted devices and subject to potential thefts. Leaked models can cause both severe financial loss and security consequences.

To understand the this new security threats, we presents the first empirical study of ML model protection on mobile devices. Our study aims to answer three open questions with quantitative evidence: How widely is model protection used in apps? How robust are existing model protection techniques? What impacts can (stolen) models incur? To that end, we built a simple app analysis pipeline and analyzed 46,753 popular apps collected from the US and Chinese app markets. We discuss our analysis and findings in detail in the following sections.

#### 4.1 Analysis Overview

We built a static-dynamic app analysis pipeline. The workflow of our analysis is depicted in Figure2. Apps first go through the static analyzer, ModelXRay, which detects the use of on-device ML and examines the model protection, if any, adopted by the app. For apps with encrypted models, the pipeline automatically generates the analysis scripts and send them to the dynamic analyzer, ModelXtractor, which performs a non-sophisticated form of in-memory extraction of model representations. ModelXtractor represents a realistic attacker who attempts to steal the decrypted ML models from an app installed on her own phone.

We report our findings and insights produced by ModelXRay and ModelXtractor in §4.2.2 and §4.3.2, respectively. We discuss the financial and security impact of (stolen) models on different stakeholders in §4.4.

#### 4.2 Q1: How Widely Is Model Protection Used in Apps?

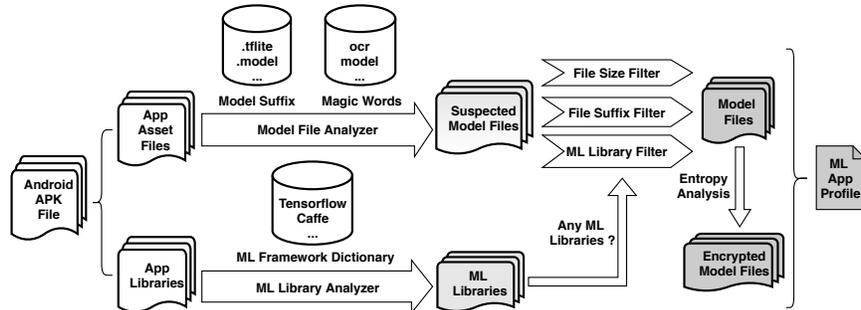


Figure 3: Identify Encrypted Models with ModelXRay

ModelXRay extracts an app’s asset files and libraries from the APK file, analyzes the native libraries and asset files to identify ML frameworks, SDK libraries and model files. Then it applies model filters combining file sizes, file suffixes and ML libraries to reduce false positives and use entropy analysis to identify encrypted models.

### 4.2.1 Methodology

We collect apps from three Android app markets: Google Play, Tencent My App, and 360 Mobile Assistant. They are the leading Android app stores in the US and China [13]. We download the apps labeled *TRENDING* and *NEW* across all 55 categories from Google Play (12,711), and all recently updated apps from Tencent My App (2,192) and 360 Mobile Assistant (31,850).

The workflow of ModelXRay is shown in Figure 3. For a given app, ModelXRay disassembles the APK file and extracts the app asset files and the native libraries. Next, it identifies the ML libraries/frameworks and the model files by analyzing the library’s exported symbols and model file’s suffix and path names. We only consider encrypted models as protected in this study. ModelXRay use entropy as an indicator of encryption. High file entropy means the file either stores encrypted data or compressed data. ModelXRay filter out compressed data by checking file magic number.

ModelXRay generates a profile for each app. A profile comprises of the information about the ML models and SDK libraries. For ML models, it records file names, sizes, MD5 hash and entropy. For SDK libraries, we record framework names, the exported symbols, and the strings extracted from the binaries. They contain information about the ML functionalities, such as OCR, face detection, liveness detection. Our analysis pipeline uses such information to generate the statistics on the use of ML libraries.

### 4.2.2 Findings and Insights

Before we present our answers to the question “Q1: How widely is model protection used in apps?”. We first show the popularity and diversity of on-device ML among our collected apps, which echo the importance of model security and protection.

**Popularity and Diversity of ML Apps:** In total, we are able to collect 46,753 Android apps from Google Play, Tencent My App and 360 Mobile Assistant stores. Using ModelXRay, we identify 1,468 apps that use on-device ML and have ML models deployed on devices, which accounts for 3.14% of our entire app collection.

*On-device ML is gaining popularity in all categories.* There are more than 50 ML apps in each of the categories, which suggests the wide interests among app developers in using on-device ML.

Table 4: The number of apps collected across markets.

Category	Google Play		Tencent My App		360 Mobile Assistant		Total	
	All	ML	All	ML	All	ML	All	ML
<b>Business</b>	404	2	99	2	2,450	296	2,953	<b>300</b>
<b>News</b>	96	0	102	5	2,450	180	2,648	<b>185</b>
<b>Images</b>	349	36	158	23	4,900	156	5,407	<b>215</b>
Map	263	4	206	14	2,450	83	2,919	101
Social	438	23	141	17	2,450	79	3,029	119
Shopping	183	5	112	16	2,450	84	2,745	105
Life	1,715	15	193	16	2,450	53	4,358	84
Education	389	3	116	7	2,450	74	2,955	84
Finance	123	6	76	21	2,450	55	2,649	82
Health	317	5	115	3	2,450	42	2,882	50
Other	8,434	79	874	35	4,900	29	14,208	143
<b>Total</b>	12,711	178	2,192	159	31,850	1,131	46,753	1,468

*Note:* In 360 Mobile Assistant, the number of unique apps is 31,591 (smaller than 32,850) because some apps are multi-categorized. Image category contains 4,900 apps because we merged image and photo related apps.

We measure the diversity of ML apps in terms of ML frameworks and functionalities. We show the top-10 most common functionalities and their distribution across different ML frameworks in Table 5.

*On-device ML offers highly diverse functionalities.* Common ML functionalities offered in the on-device fashion includes OCR, face tracking, hand detection, speech recognition, handwriting recognition, ID card recognition, and bank card recognition, liveness detection, face recognition, iris recognition and so on.

**Model Protection Across App Stores:** Figure 4 gives the per-app-market statistics on ML model protection and reuse. Figure 4a shows the per-market numbers of protected apps (*i.e.*, apps using protected/encrypted models)

Table 5: Number of apps using different ML Frameworks with different functionalities.

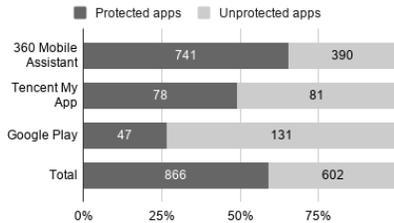
Functionality	TensorFlow (Google)	*Caffe2/PyTorch (Facebook)	*Parrots (SenseTime)	TFLite (Google)	NCNN (Tencent)	Mace (Xiaomi)	MxNet (Apache)	ULS (Utility Asset Store)	Total
OCR(Optical Character Recognition)	41	186	140	6	37	18	1	11	441
Face Tracking	26	272	216	7	53	6	13	27	620
Speech Recognition	7	32	9	1	11	18	1	9	88
Hand Detection	4	0	0	2	4	0	0	0	10
Handwriting Recognition	8	17	1	0	16	0	0	0	42
<b>Liveness Detection</b>	32	392	349	9	70	7	10	3	872
<b>Face Recognition</b>	17	116	95	6	40	7	10	3	294
<b>Iris Recognition</b>	0	4	0	0	2	0	3	0	9
ID Card Recognition	26	230	147	5	47	18	0	10	483
Bank Card Recognition	11	126	117	2	16	18	0	9	299

Note: 1) One app may use multiple frameworks for different ML functionalities. Therefore, the sum of apps using different functionalities is bigger than the number of total apps. 2) Security critical functionalities are in **bold fonts** and can be used for fraud detection or access control. 3) \*Caffe was initially developed by Berkeley, based on which Facebook built Caffe2, which was later merged with PyTorch. The following uses “Caffe” to represent Caffe, Caffe2 and PyTorch.

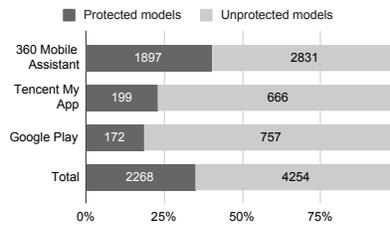
and unprotected apps (*i.e.*, apps using unprotected models).

Overall, only 59% of ML apps protect their models. The rest of the apps (602 in total) simply include the models in plaintext, which can be easily extracted from the app packages or installation directories. For 41% of the ML apps, stealing their models is as easy as downloading and decompressing their app packages.

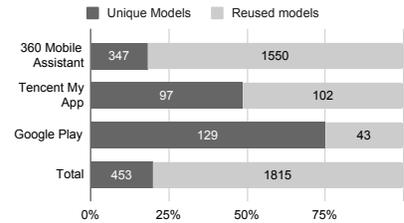
If we focus on individual models (*i.e.*, some apps use multiple ML models for different functionalities), the percentages of unprotected models (Figure 4b) become even higher. Overall, 4,254 out of 6,522 models (77%) are unprotected and thus easily extractable and reverse engineered.



(a) Apps using protected/encrypted models vs. those using unprotected models

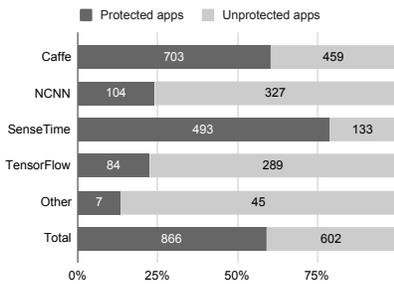


(b) On-device models that are protected/encrypted vs. those not

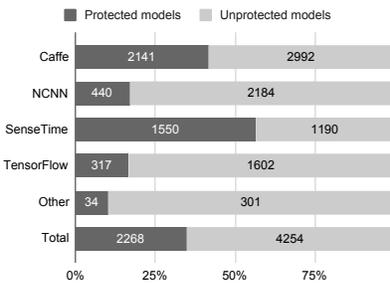


(c) Unique encrypted models vs. encrypted models reused/shared by multiple apps.

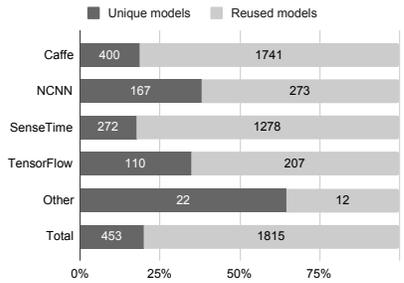
Figure 4: Statistics on ML model protection and reuse, grouped by app markets. The “total” number of unique models is less than the sum of the per-store numbers because some models are not unique from different stores.



(a) Apps using protected models vs. those using unprotected models



(b) On-device models that are protected/encrypted vs. those not



(c) Unique encrypted models vs. encrypted models reused/shared by multiple apps

Figure 5: Statistics on ML model protection and reuse, grouped by ML frameworks. The “total” number is less than the sum of the per-framework numbers because many apps use multiple frameworks for different functionalities.

**Model Protection Across ML Frameworks:** Some ML frameworks have wider adoption of model protection.

Figure 5 shows that, more than 79% of the apps using SenseTime (Parrots) have protected models, followed by

apps using Caffe (60% of them have protected models). For apps using TensorFlow and NCNN, the number is around 20%. Apps using other frameworks are the least protected against model thefts.

**GPU Acceleration Adoption Rate among ML Apps:** Table 6 shows the number ML apps and libraries that use GPU for acceleration. 797(54%) ML apps make use of GPU. *The wide adoption of GPU acceleration poses a challenge to the design of secure on-device ML.* For instance, the naive idea of performing model inference and other model access operations entirely inside a trusted execution environment (TEE, e.g., TrustZone) is not viable due to the need for GPU acceleration, which cannot be easily or efficiently accessed within the TEE.

Table 6: ML apps and libraries that use GPU acceleration

	360 Mobile Assistant	Tencent My App	Google Play
ML Apps	669	104	24
ML Libraries	212	103	23

### 4.3 Q2: How Robust Are Existing Model Protection Techniques?

We build ModelXtractor, a tool simple by design to dynamically recover protected or encrypted models used in on-device ML. ModelXtractor mainly targets on-device ML models that are encrypted during transportation and at rest (in storage) but not protected when in use or loaded in memory. For protected models mentioned in §4.2, ModelXtractor is performed to assess the robustness of the protection.

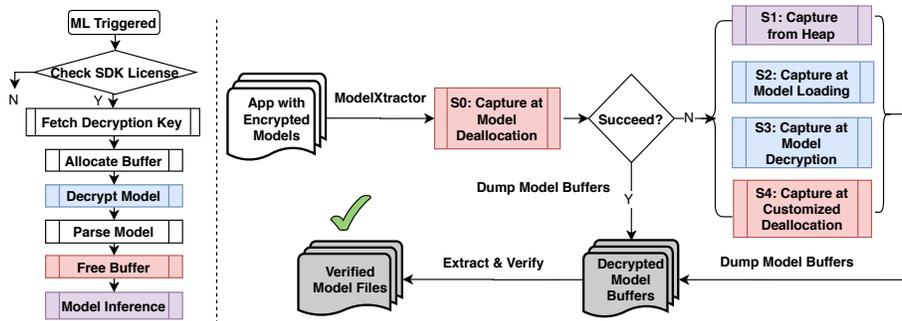


Figure 6: Extraction of (decrypted) models from app memory using ModelXtractor

The left side shows the typical workflow of model loading and decryption in mobile apps. The right side shows the workflow of ModelXtractor. The same color on both sides indicate the same timing of the strategy being used. The "Check SDK License" shows that a model provider will check an app's SDK license before releasing the decryption keys as a way to protect its IP.

The workflow of ModelXtractor is depicted in Figure 6. It takes inputs from ModelXRay, including the information about the ML framework(s) and the model(s) used in the app (described in §4.2). These information helps to target and efficiently instrument an app during runtime, and capture models in plaintext from the memory of the app.

#### 4.3.1 Methodology

ModelXtractor uses app instrumentation to dynamically find the memory buffers where (decrypted) ML is loaded and accessed by the ML frameworks. For each app, ModelXtractor determines which libraries and functions need to be instrumented and when to start and stop each instrumentation, based on the instrumentation strategies (discussed shortly). ModelXtractor automatically generates the code that needs to be inserted at different instrumentation points. It employs the widely used Android instrumentation tool, Frida [6], to perform code injection.

ModelXtractor has a main instrumentation strategy (S0) and four alternative ones (S1-S4). When the default strategy cannot capture the models, the alternatively strategies (S1-S4) will be used.

**S0: Capture at Model Deallocation:** This is the default strategy since we observe the most convenient time and place to capture an in-memory model is right before the deallocation of the buffer where the model is loaded.

This is because (1) memory deallocation APIs (*e.g.*, `free`) are limited in numbers and easy to instrument, and (2) models are completely loaded and decrypted when their buffers are to be freed.

This default instrumentation strategy may fail in the following uncommon scenarios. First, an app is not using native ML libraries, but a JavaScript ML library. Second, an app uses its own or customized memory allocator/deallocator. Third, a model buffer is not freed during our dynamic analysis.

**S1: Capture from Heap:** This strategy dumps the entire heap region of an app when a ML functionality is in use, in order to identify possible models in it. It is suitable for apps that do not free model buffers timely or at all. It also helps in cases where memory-managed ML libraries are used (*e.g.*, JavaScript) and buffer memory deallocations (done by a garbage collector) are implicit or delayed.

**S2: Capture at Model Loading:** This strategy instruments ML framework APIs that load models to buffers. We manually collect a list of such APIs (*e.g.*, `loadModel`) for the ML frameworks observed in our analysis. This strategy is suitable for those apps where S0 fails and the ML framework code is not obfuscated.

**S3: Capture at Model Decryption:** This strategy instruments model decryption APIs (*e.g.*, `aes256_decrypt`) in ML frameworks, which we collected manually. Similar to S2, it is not applicable to apps that use obfuscated ML framework code.

**S4: Capture at Customized Deallocation:** Some apps use customized memory deallocators. We manually identify a few such allocators (*e.g.*, `slab_free`), which are instrumented similarly as S0.

### 4.3.2 Findings and Insights

**Results of Dynamic Model Extraction:** Table 7 shows the statistics on the 82 analyzed apps, grouped by the ML frameworks they use. Among the 29 apps whose ML functionalities were triggered, we successfully extracted models from 18 of them (66%). Considering the reuse of those extracted encrypted models, the number of apps that are affected by our model extraction is 347 (*i.e.*, 347 apps used the same models and same protection techniques as the 18 apps that we extracted models from). This extraction rate is alarming and shows that a majority of the apps using model protection can still lose their valuable models to an unsophisticated attack. It indicates that *even for app developers and ML providers willing/trying to protect their models, it is hard to do it in a robust way using the file encryption-based techniques.*

Table 8 shows the per-app details about the extracted models. We anonymized the apps for security concerns: many of them are highly downloaded apps or provide security-critical services. Many of the listed apps contain more than one ML models. For simplicity, we only list one representative model for each app.

*The extracted models are highly popular and diverse, some very valuable or security-critical.* From Table 8 we can see that 8 of 15 listed apps have been downloaded more than 10 million times. Half of the extracted models belong to commercial ML providers, such as SenseTime, and were purchased by the app developers. Such models being leaked may cause direct financial loss to both app developers and model owners (§4.4).

As for diversity, the model size ranges from 160KB to 20MB. They span all the popular frameworks, such as TensorFlow, TFLite, Caffe, SenseTime, Baidu, and Face++. The observed model formats include Protobuf, FlatBuffer, JSON, and some proprietary formats used by SenseTime, Face++ and Baidu. In terms of ML functionalities, the models are used for face recognition, face tracking, liveness detection, OCR, ID/card recognition, photo processing, and malware detection. Among them, liveness detection, malware detection, and face recognition are often used for security-critical purposes, such as access control and fraud detection. Leakage of these models may give attackers an advantage to develop model evasion techniques in a white-box fashion.

Table 7: Model extraction statistics.

ML Framework	Unique Models Analyzed	ML Triggered	Models Extracted	Models Missed	Apps Affected
TensorFlow	3	3	3	0	3
Caffe	7	3	1	2	79
SenseTime	55	16	11	5	186
TFLite	3	2	2	0	76
NCNN	9	3	0	3	0
Other	5	3	2	1	88
<b>Total</b>	<b>82</b>	<b>29</b>	<b>18</b>	<b>11</b>	<b>347</b>

*Note:* 347 is the sum of affected apps per framework after deduplication.

Table 8: Overview of Successfully Dumped Models with ModelXtractor

App name	Downloads	Framework	Model Functionality	Size (B)	Format	Reuses	Extraction Strategy
Anonymous App 1	300M	TFLite	Liveness Detection	160K	FlatBuffer	18	Freed Buffer
Anonymous App 2	10M	Caffe	Face Tracking	1.5M	Protobuf	4	Model Loading
Anonymous App 3	27M	SenseTime	Face Tracking	2.3M	Protobuf	77	Freed Buffer
Anonymous App 4	100K	SenseTime	Face Filter	3.6M	Protobuf	3	Freed Buffer
Anonymous App 5	100M	SenseTime	Face Filter	1.4M	Protobuf	2	Freed Buffer
Anonymous App 6	10K	TensorFlow	OCR	892K	Protobuf	2	Memory Dumping
Anonymous App 7	10M	TensorFlow	Photo Process	6.5M	Protobuf	1	Freed Buffer
Anonymous App 8	10K	SenseTime	Face Track	1.2M	Protobuf	5	Freed Buffer
Anonymous App 9	5.8M	Caffe	Face Detect	60K	Protobuf	77	Freed Buffer
Anonymous App 10	10M	Face++	Liveness	468K	Unknown	17	Freed Buffer
Anonymous App 11	100M	SenseTime	Face Detect	1.7M	Protobuf	18	Freed Buffer
Anonymous App 12	492K	Baidu	Face Tracking	2.7M	Unknown	26	Freed Buffer
Anonymous App 13	250K	SenseTime	ID card	1.3M	Unknown	13	Freed Buffer
Anonymous App 14	100M	TFLite	Camera Filter	228K	Json	1	Freed Buffer
Anonymous App 15	5K	TensorFlow	Malware Classification	20M	Protobuf	1	Decryption Buffer

Note: 1) We excluded some apps that dumped the same models as reported above; 2) We anonymized the name of the apps to protect the user’s security; 3) Every app has several models for different functionalities, we only list one representative model for each app.

## 4.4 Q3: What Impacts can (Stolen) Models Incur?

ML models are the core intellectual properties of ML solution providers. The impacts of leaked models are wide and profound, including substantial financial impact as well as significant security implications.

### 4.4.1 Financial Impact

Financial impact mainly applies to two stakeholders: attackers, and vendors (model providers and app companies).

**Attackers financially benefit from leaked models.** App developers usually have two legitimate ways to get ML models: (1) buying a ML SDK and model license from a ML solution provider, such as SenseTime, Face++, and so on; (2) designing and training their own ML models using open-source or customized frameworks, which usually requires a large amount of computing and human resources. Stealing the models saves the attackers either the license fee paid to the model providers, or the research and development (R&D) cost on the models.

According to Face++, the annual fee for a license with offline authorization is \$50,000 to \$200,000 [8]. The saving is large enough to motivate an attacker to steal the models. In our analysis, we found 60 cases of different app companies are reusing model licenses. One of the licenses is even used by potentially 12 different app companies, indicating a high chance of illegal uses.

**Vendors face financial loss from decreasing competitive edge with stolen models.** For vendors whose main business (source of income) depends on ML models, e.g., model providers or app companies, model leakages result into pricing disadvantages, lost of customers and market share.

For model providers, the market is strongly competitive. In our study, we have found some top ML SDK providers, such as SenseTime, Megvii, Baidu, ULSee, Anyline, etc. Take Megvii as an example, according to Owler [9], 10 competitors are closely related to its businesses, such as Cognitec, SenseTime, Kairos, FaceFirst, Cortexica, etc. For app companies, the competition is as much competitive if not more so. In Google Play only, our study found 36 apps using ML SDK for image recognition as the main business. Considering the other two stores, at least 215 apps are competing for this business.

### 4.4.2 Security Impact

Some ML models are used for security-critical purposes. For example, liveness detection model is used to verify whether it is a real person holding a real ID card. Face, fingerprint and iris recognition models are used to detect and verify the identity of a person. These models bring in great convenience, for example, users do not need to go to a bank or customer service centers to verify their identities. However, breaches of such models bring in security and privacy concerns.

For attackers, a leaked security-critical model makes it easier for them to design and craft adversarial examples. They can then use the examples to either fake different identities, or simply bypass the identity check of the apps [3].

We found more than 100 apps using on-device ML models for banking and loan services. These apps provide personal loan services aiming at quick and convenient loan applications. They use face recognition models to verify

the identity of a person by taking a short video, and comparing with the photo on the ID card. The apps then determine the credit limits and rates to loan to the applicants. When the models are leaked, attackers can easily fake identities of other applicants, and apply for loans on their behalf.

## 5 Research Plan

In previous sections, I presented the completed work on defending advanced attacks on IoT devices and understanding the newly emerging model privacy problem on mobile devices. However, how to protect the model privacy on mobile devices is still an open research problem.

Our previous work [67] on model privacy shows that even though the ML models can be encrypted by the app, it still suffers from unsophisticated runtime attacks. 54% ML apps use GPU acceleration, which means supporting GPU acceleration is very important for existing ML apps. Besides, many ML apps use on-device ML for live video stream analysis including face recognition, liveness detection and so on. With all above consideration, we want to design a secure on-device model inference system that meets the following goals:

- Security rooted in the hardware, so that the models will be secure even when the OS is compromised;
- Reasonable performance overhead, making sure it will not break existing real-time tasks;
- Access to hardware accelerators, which are being developed and used for on-device ML tasks;

TEE provides hardware-level security. Our previous work OAT [66] demonstrated that TEE can be used to build efficient attestation system with hardware-backed trusted computing primitives. However, using mobile TEE for secure model inference has several technical challenges.

First, the mobile TEE, like Arm TrustZone, is designed for small security critical services like managing encryption keys. The memory reserved for the TEE OS is limited. For example, only 14 MB is available for Trusted Applications of OP-TEE OS on Hikey960 Dev Board, while the model size of AlexNet is 242MB. It is simply not feasible to run the high resource-demanding model inference inside the TEE.

Second, the current TEE does not include GPU/NPU into the secure domain, so we will also lose access to the hardware acceleration. Third, the model inference framework will also greatly increase the TCB of TEE, risking the security of the whole system.

To address the above challenges, I propose the following work as part of my doctoral thesis.

### 5.1 ShadowNet: a Secure and Efficient Model Inference Scheme

We propose ShadowNet: a secure and efficient model inference scheme built on top of mobile TEE. The key observation of ShadowNet is shared with the previous research like Slalom [68] that the linear layers of CNNs occupy the majority of the model parameters and the model inference time. For example, the linear layers of MobileNets occupy around 95% of the model parameters, 99% of the model inference time. The idea of ShadowNet is to apply linear transformation on the weights of the linear layers and outsource them onto the untrusted world, so that we can leverage the hardware acceleration without trusting it. ShadowNet restores the results inside the TEE. The other nonlinear layers are also kept secure inside the TEE.

With this design, ShadowNet’s security is rooted in the TEE, meeting the first design goal; ShadowNet does not introduce any heavy cryptographic operations, the performance overhead should not be heavy, meeting the second design goal; ShadowNet is still able to use the hardware acceleration meeting the third design goal. At the same time, ShadowNet solves the technical challenges of mobile TEE by keeping the TCB size small and the TEE memory usage low.

### 5.2 Formal Security Analysis of the Scheme

In order to apply ShadowNet scheme, we need to formally prove the scheme is secure. In detail, we need to show that ShadowNet meet two security properties: (1) by observing the weights of the transformed linear layers that is outsourced to the untrusted world, the attacker won’t be able to infer the original secret weights before transformation. (2) based on the exposed transformed weights, training an equivalent model should not be made significantly easier than training the original model from scratch.

These two security property guarantee that after applying the ShadowNet transformation, the transformed model will not leak the original weights or make it easier for the attacker to train an equivalent model by outsourcing the transformed linear layers onto the untrusted hardware accelerators.

### 5.3 System Design Challenges

ShadowNet split the model inference between the Rich Execution Environment (REE) and the Trusted Execution Environment (TEE). The transformed linear layers can run on the untrusted hardware accelerators including GPU while the restoration of the transformed results and the model inference of non-linear layers run inside the TEE.

In order to achieve good performance, we need to overcome several system design and engineering challenges. First, the communication overhead between REE and TEE should be optimized as it happens frequently. The parameter passing between REE and TEE should be designed in a secure and efficient way. Second, writing code for TEE also faces the challenges of limited resources, like limited secure memory. The model inference code should be written in a memory efficient way. Third, mobile TEE does not have rich computation library support. For example, several popular computation libraries including Eigen Library [7] and Arm Compute Library [2] only have C++ version and do not support mobile TEE OS like the OP-TEE OS [56].

### 5.4 Milestones

The plan for completing my research is presented in Table 9.

Table 9: Plan for completion

<b>Task</b>	<b>Completion Date</b>
Finish implementation & evaluation of ShadowNet	January 2021
Formalize the security analysis of ShadowNet	February 2021
Finalizing ShadowNet for publication	March 2021
Dissertation defense	April 2021

## References

- [1] A brief guide to mobile AI chips. <https://www.theverge.com/2017/10/19/16502538/mobile-ai-chips-apple-google-huawei-qualcomm>.
- [2] Arm Compute Library. <https://www.arm.com/why-arm/technologies/compute-library>.
- [3] Artificial Intelligence + GANs can create fake celebrity faces. <https://medium.com/datadriveninvestor/artificial-intelligence-gans-can-create-fake-celebrity-faces-44fe80d419f7>.
- [4] Asylo - An open and flexible framework for enclave applications. <https://asylo.dev/>.
- [5] Converting model to C++ code. [https://mace.readthedocs.io/en/latest/user\\_guide/advanced\\_usage.html](https://mace.readthedocs.io/en/latest/user_guide/advanced_usage.html).
- [6] Dynamic instrumentation toolkit for developers, reverse-engineers, and security researchers. <https://frida.re/>.
- [7] Eigen. [http://eigen.tuxfamily.org/index.php?title=Main\\_Page](http://eigen.tuxfamily.org/index.php?title=Main_Page).
- [8] Face++ pricing details - mobile sdk. <https://www.faceplusplus.com/pricing-details/#offline>.
- [9] Megvii's Competitors, Revenue, Number of Employees, Funding and Acquisitions. <https://www.owler.com/company/megvii>.
- [10] Strip visible string in ncnn. <https://github.com/Tencent/ncnn/wiki>.
- [11] TF Encrypted. <https://github.com/tf-encrypted/tf-encrypted>.
- [12] TF Trusted. <https://github.com/dropoutlabs/tf-trusted>.
- [13] The AppInChina App Store Index. <https://www.appinchina.co/market/app-stores/>, 2019.
- [14] Tigist Abera, N. Asokan, Lucas Davi, Jan-Erik Ekberg, Thomas Nyman, Andrew Paverd, Ahmad-Reza Sadeghi, and Gene Tsudik. C-flat: Control-flow attestation for embedded systems software. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, CCS '16, pages 743–754, New York, NY, USA, 2016. ACM.
- [15] Periklis Akrividis, Cristian Cadar, Costin Raiciu, Manuel Costa, and Miguel Castro. Preventing memory error exploits with wit. In *Security and Privacy, 2008. SP 2008. IEEE Symposium on*, pages 263–277. IEEE, 2008.
- [16] W. A. Arbaugh, D. J. Farber, and J. M. Smith. A secure and reliable bootstrap architecture. In *Proceedings of the 1997 IEEE Symposium on Security and Privacy*, SP '97, pages 65–, Washington, DC, USA, 1997. IEEE Computer Society.
- [17] Frederik Armknecht, Ahmad-Reza Sadeghi, Steffen Schulz, and Christian Wachsmann. A security framework for the analysis and design of software attestation. In *Proceedings of the 2013 ACM SIGSAC Conference on Computer & Communications Security*, CCS '13, pages 1–12, New York, NY, USA, 2013. ACM.
- [18] Sergei Arnautov, Bohdan Trach, Franz Gregor, Thomas Knauth, Andre Martin, Christian Priebe, Joshua Lind, Divya Muthukumar, Dan O'keeffe, Mark L Stillwell, et al. {SCONE}: Secure linux containers with intel {SGX}. In *12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16)*, pages 689–703, 2016.
- [19] N. Asokan, Ferdinand Brasser, Ahmad Ibrahim, Ahmad-Reza Sadeghi, Matthias Schunter, Gene Tsudik, and Christian Wachsmann. Seda: Scalable embedded device attestation. In *Proceedings of the 22Nd ACM SIGSAC Conference on Computer and Communications Security*, CCS '15, pages 964–975, New York, NY, USA, 2015. ACM.
- [20] Sebastian P Bayerl, Tommaso Frassetto, Patrick Jauernig, Korbinian Riedhammer, Ahmad-Reza Sadeghi, Thomas Schneider, Emmanuel Stapf, and Christian Weinert. Offline model guard: Secure and private ml on mobile devices. *DATE 2020*, 2020.
- [21] Ferdinand Brasser, David Gens, Patrick Jauernig, Ahmad-Reza Sadeghi, and Emmanuel Stapf. Sanctuary: Arming trustzone with user-space enclaves. In *NDSS*, 2019.
- [22] Erik Buchanan, Ryan Roemer, Hovav Shacham, and Stefan Savage. When good instructions go bad: Generalizing return-oriented programming to risc. In *Proceedings of the 15th ACM conference on Computer and communications security*, pages 27–38. ACM, 2008.
- [23] Scott A Carr and Mathias Payer. Datashield: Configurable data confidentiality and integrity. In *Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security*, pages 193–204. ACM, 2017.
- [24] Claude Castelluccia, Aurélien Francillon, Daniele Perito, and Claudio Soriente. On the difficulty of software-based attestation of embedded devices. In *Proceedings of the 16th ACM Conference on Computer and Communications Security*, CCS '09, pages 400–409, New York, NY, USA, 2009. ACM.
- [25] Miguel Castro, Manuel Costa, and Tim Harris. Securing software by enforcing data-flow integrity. In *Proceedings of the 7th Symposium on Operating Systems Design and Implementation*, OSDI '06, pages 147–160, Berkeley, CA, USA, 2006. USENIX Association.
- [26] Huili Chen, Cheng Fu, Bitar Darvish Rouhani, Jishen Zhao, and Farinaz Koushanfar. DeepAttest: An End-to-End Attestation Framework for Deep Neural Networks. 2019.
- [27] Shuo Chen, Jun Xu, Emre Can Sezer, Prachi Gauriar, and Ravishankar K Iyer. Non-control-data attacks are realistic threats. In *USENIX Security Symposium*, volume 14, 2005.
- [28] Xuhui Chen, Jinlong Ji, Lixing Yu, Changqing Luo, and Pan Li. Securenets: Secure inference of deep neural networks on an untrusted cloud. In *Asian Conference on Machine Learning*, pages 646–661, 2018.
- [29] Karim Eldefrawy, Aurélien Francillon, Daniele Perito, and Gene Tsudik. SMART: Secure and Minimal Architecture for (Establishing a Dynamic) Root of Trust. In *NDSS 2012, 19th Annual Network and Distributed System Security Symposium, February 5-8, San Diego, USA*, San Diego, UNITED STATES, 02 2012.
- [30] Nicolas Falliere, Liam O Murchu, and Eric Chien. W32. stuxnet dossier. *White paper, Symantec Corp., Security Response*, 5(6), 2011.
- [31] FDA. Cybersecurity vulnerabilities identified in implantable cardiac pacemaker, August 2017.

- [32] Zahra Ghodsi, Tianyu Gu, and Siddharth Garg. Safetynets: Verifiable execution of deep neural networks on an untrusted cloud. In *Advances in Neural Information Processing Systems*, pages 4672–4681, 2017.
- [33] Dennis Giese. Security analysis of the xiaomi iot ecosystem. Master’s thesis, Master’s thesis, Technische Universität Darmstadt, 2019.
- [34] Ran Gilad-Bachrach, Nathan Dowlin, Kim Laine, Kristin Lauter, Michael Naehrig, and John Wernsing. Cryptonets: Applying neural networks to encrypted data with high throughput and accuracy. In *International Conference on Machine Learning*, pages 201–210, 2016.
- [35] Github. House Alarm System. <https://github.com/ddrazir/alarm4pi>.
- [36] Github. Light Controller. <https://github.com/Barro/light-controller>.
- [37] Github. Remote Movement Controller. <https://github.com/bskari/pi-rc/tree/pi2>.
- [38] Github. Rover Controller. <http://github.com/Gwaltrip/RoverPi/tree/master/tcpRover>.
- [39] Andy Greenberg. Hacker says he can hijack a \$35 k police drone a mile away, 2016.
- [40] Zhongshu Gu, Heqing Huang, Jialong Zhang, Dong Su, Ankita Lamba, Dimitrios Pendarakis, and Ian Molloy. Yerbabuena: Securing deep learning inference data via enclave-based ternary model partitioning. *arXiv preprint arXiv:1807.00969*, 2018.
- [41] Lucjan Hanzlik, Yang Zhang, Kathrin Grosse, Ahmed Salem, Max Augustin, Michael Backes, and Mario Fritz. Mlcapsule: Guarded offline deployment of machine learning as a service. *arXiv preprint arXiv:1808.00590*, 2018.
- [42] Grant Ho, Derek Leung, Pratyush Mishra, Ashkan Hosseini, Dawn Song, and David Wagner. Smart locks: Lessons for securing commodity internet of things devices. In *Proceedings of the 11th ACM on Asia Conference on Computer and Communications Security*, pages 461–472. ACM, 2016.
- [43] Hong Hu, Shweta Shinde, Sendriou Adrian, Zheng Leong Chua, Prateek Saxena, and Zhenkai Liang. Data-oriented programming: On the expressiveness of non-control data attacks. In *Security and Privacy (SP), 2016 IEEE Symposium on*, pages 969–986. IEEE, 2016.
- [44] Troy Hunt. Controlling vehicle features of nissan leafs across the globe via vulnerable apis. <https://www.troyhunt.com/controlling-vehicle-features-of-nissan/>, February 2016.
- [45] Ahmad Ibrahim, Ahmad-Reza Sadeghi, Gene Tsudik, and Shaza Zeitouni. Darpa: Device attestation resilient to physical attacks. In *Proceedings of the 9th ACM Conference on Security & Privacy in Wireless and Mobile Networks*, WiSec ’16, pages 171–182, New York, NY, USA, 2016. ACM.
- [46] Xiaoqian Jiang, Miran Kim, Kristin Lauter, and Yongsoo Song. Secure outsourced matrix computation and application to neural networks. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*, pages 1209–1222, 2018.
- [47] Chongkyung Kil, Emre C. Sezer, Ahmed M. Azab, Peng Ning, and Xiaolan Zhang. Remote attestation to dynamic system properties: Towards providing complete system integrity evidence. *Dependable Systems & Networks, 2009.*, 2009.
- [48] Joonho Kong, Farinaz Koushanfar, Praveen K. Pendyala, Ahmad-Reza Sadeghi, and Christian Wachsmann. PUFatt: Embedded Platform Attestation Based on Novel Processor-Based PUFs. In *DAC*, page 6. ACM, 2014.
- [49] Tim Kornau. *Return oriented programming for the ARM architecture*. PhD thesis, Master’s thesis, Ruhr-Universität Bochum, 2010.
- [50] Brian Krebs. Reaper: Calm before the iot security storm. <https://krebsonsecurity.com/2017/10/reaper-calm-before-the-iot-security-storm/>, October 2017.
- [51] Roland Kunkel, Do Le Quoc, Franz Gregor, Sergei Arnautov, Pramod Bhatotia, and Christof Fetzer. TensorSCONE: A Secure TensorFlow Framework using Intel SGX. *arXiv preprint arXiv:1902.04413*, 2019.
- [52] Keen Security Lab. New car hacking research: Tesla motors. <http://keenlab.tencent.com/en/2017/07/27/New-Car-Hacking-Research-2017-Remote-Attack-Tesla-Motors-Again/>, 2017.
- [53] Taegyong Lee, Zhiqi Lin, Saumay Pushp, Caihua Li, Yunxin Liu, Youngki Lee, Fengyuan Xu, Chenren Xu, Lintao Zhang, and Junehwa Song. Occlumency: Privacy-preserving remote deep-learning inference using sgx. In *The 25th Annual International Conference on Mobile Computing and Networking*, pages 1–17, 2019.
- [54] Matthew Lentz, Rijurekha Sen, Peter Druschel, and Bobby Bhattacharjee. Secloak: Arm trustzone-based mobile peripheral control. In *Proceedings of the 16th Annual International Conference on Mobile Systems, Applications, and Services*, pages 1–13, 2018.
- [55] Yanlin Li, Jonathan M. McCune, and Adrian Perrig. VIPER: verifying the integrity of peripherals’ firmware. In *Proceedings of the 18th ACM Conference on Computer and Communications Security, CCS 2011, Chicago, Illinois, USA, October 17-21, 2011*, pages 3–16, 2011.
- [56] Linaro. OP-TEE. <https://www.op-tee.org>.
- [57] Fan Mo, Ali Shahin Shamsabadi, Kleomenis Katevas, Soteris Demetriou, Ilias Leontiadis, Andrea Cavallaro, and Hamed Haddadi. Darknetz: Towards model privacy at the edge using trusted execution environments. In *ACM MobiSys 2020*.
- [58] Bryan Parno, Jonathan M. McCune, and Adrian Perrig. Bootstrapping trust in commodity computers. In *Proceedings of the 2010 IEEE Symposium on Security and Privacy, SP ’10*, pages 414–429, Washington, DC, USA, 2010. IEEE Computer Society.
- [59] D. Quarta, M. Pogliani, M. Polino, F. Maggi, A. M. Zanchettin, and S. Zanero. An experimental security analysis of an industrial robot controller. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 268–286, May 2017.
- [60] Rapid7. Multiple vulnerabilities in animas onetouch ping insulin pump. <https://blog.rapid7.com/2016/10/04/r7-2016-07-multiple-vulnerabilities-in-animas-onetouch-ping-insulin-pump/>, October 2016.

- [61] Eyal Ronen, Adi Shamir, Achi-Or Weingarten, and Colin O’Flynn. Iot goes nuclear: Creating a zigbee chain reaction. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 195–212. IEEE, 2017.
- [62] Sergio Salinas, Changqing Luo, Weixian Liao, and Pan Li. Efficient secure outsourcing of large-scale quadratic programs. In *Proceedings of the 11th ACM on Asia Conference on Computer and Communications Security*, pages 281–292, 2016.
- [63] A. Seshadri, A. Perrig, L. van Doorn, and P. Khosla. Swatt: software-based attestation for embedded devices. In *IEEE Symposium on Security and Privacy, 2004. Proceedings. 2004*, pages 272–282, May 2004.
- [64] Arvind Seshadri, Mark Luk, Elaine Shi, Adrian Perrig, Leendert van Doorn, and Pradeep Khosla. Pioneer: Verifying code integrity and enforcing untampered code execution on legacy systems. *SIGOPS Oper. Syst. Rev.*, 39(5):1–16, October 2005.
- [65] Christos Stergiou, Kostas E Psannis, Byung-Gyu Kim, and Brij Gupta. Secure integration of iot and cloud computing. *Future Generation Computer Systems*, 78:964–975, 2018.
- [66] Zhichuang Sun, Bo Feng, Long Lu, and Somesh Jha. Oat: Attesting operation integrity of embedded devices. In *2020 IEEE Symposium on Security and Privacy (SP)*, pages 1433–1449. IEEE, 2020.
- [67] Zhichuang Sun, Ruimin Sun, and Long Lu. Mind your weight (s): A large-scale study on insufficient machine learning model protection in mobile apps. *arXiv preprint arXiv:2002.07687*, 2020.
- [68] Florian Tramer and Dan Boneh. Slalom: Fast, verifiable and private execution of neural networks in trusted hardware. *arXiv preprint arXiv:1806.03287*, 2018.
- [69] Yusuke Uchida, Yuki Nagai, Shigeyuki Sakazawa, and Shin’ichi Satoh. Embedding watermarks into deep neural networks. In *Proceedings of the 2017 ACM on International Conference on Multimedia Retrieval*, pages 269–277. ACM, 2017.
- [70] Arvind Seshadri Adrian Perrig Leendert van Doorn Pradeep Khosla. Using software-based attestation for verifying embedded systems in cars. *S&P, Oakland*, 2004.
- [71] Deepak Chandra Vivek Haldar and Michael Franz. Semantic remote attestation – a virtual machine directed approach to trusted computing. *VM*, 2004.
- [72] Stavros Volos, Kapil Vaswani, and Rodrigo Bruno. Graviton: Trusted execution environments on GPUs. In *13th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 18)*, pages 681–696, 2018.
- [73] Wikipedia. Mirai(malware). [https://en.wikipedia.org/wiki/Mirai\\_\(malware\)](https://en.wikipedia.org/wiki/Mirai_(malware)), 2020.
- [74] Jacob Wurm, Khoa Hoang, Orlando Arias, Ahmad-Reza Sadeghi, and Yier Jin. Security analysis on consumer and industrial iot devices. In *Design Automation Conference (ASP-DAC), 2016 21st Asia and South Pacific*, pages 519–524. IEEE, 2016.
- [75] Jialong Zhang, Zhongshu Gu, Jiyong Jang, Hui Wu, Marc Ph Stoecklin, Heqing Huang, and Ian Molloy. Protecting intellectual property of deep neural networks with watermarking. In *Proceedings of the 2018 on Asia Conference on Computer and Communications Security*, pages 159–172. ACM, 2018.