



The Threat of Offensive AI to Organizations

Yisroel Mirsky^{a,*}, Ambra Demontis^b, Jaidip Kotak^a, Ram Shankar^c, Deng Gelei^d, Liu Yang^d, Xiangyu Zhang^e, Maura Pintor^f, Wenke Lee^g, Yuval Elovici^a, Battista Biggio^f

^a Ben-Gurion University of the Negev, Israel

^b University of Cagliari, Italy

^c Microsoft, United States of America

^d Nanyang Technological University, Singapore

^e Purdue University, United States of America

^f University of Cagliari & Pluribus One, Italy

^g Georgia Institute of Technology, United States of America

ARTICLE INFO

Article history:

Received 10 July 2022

Revised 7 October 2022

Accepted 3 November 2022

Available online 6 November 2022

Keywords:

Offensive AI

APT

Cyber security

Organization security

Adversarial machine learning

Deepfake

AI-Capable adversary

ABSTRACT

AI has provided us with the ability to automate tasks, extract information from vast amounts of data, and synthesize media that is nearly indistinguishable from the real thing. However, positive tools can also be used for negative purposes. In particular, cyber adversaries can use AI to enhance their attacks and expand their campaigns.

Although offensive AI has been discussed in the past, there is a need to analyze and understand the threat in the context of organizations. For example, how does an AI-capable adversary impact the cyber kill chain? Does AI benefit the attacker more than the defender? What are the most significant AI threats facing organizations today and what will be their impact on the future?

In this study, we explore the threat of offensive AI on organizations. First, we present the background and discuss how AI changes the adversary's methods, strategies, goals, and overall attack model. Then, through a literature review, we identify 32 offensive AI capabilities which adversaries can use to enhance their attacks. Finally, through a panel survey spanning industry, government and academia, we rank the AI threats and provide insights on the adversaries.

© 2022 Elsevier Ltd. All rights reserved.

1. Introduction

For decades, organizations, including government agencies, hospitals, and financial institutions, have been the target of cyber attacks (Knake, 2017; Mattei, 2017; Tariq, 2018). These cyber attacks have been carried out by experienced hackers using manual methods. In recent years there has been a boom in the development of artificial intelligence (AI), which has enabled the creation of software tools that have helped to automate tasks such as prediction, information retrieval, and media synthesis. Throughout this period, members of academia and industry have utilized AI¹ in the context

of improving the state of cyber defense (Liu and Lang, 2019; Mahadi et al., 2018; Mirsky et al., 2018) and threat analysis (Cohen et al., 2020; Deepreflect, 2021; Ucci et al., 2019). However, AI is a double edged sword, and attackers can utilize it to improve their malicious campaigns.

Therefore, we define Offensive AI as

“The use or abuse of AI to accomplish a malicious task”

Offensive Use of AI. Adversaries can improve their tactics to launch attacks that were not possible before. For example, with deep learning one can perform highly effective spear phishing attacks by impersonating their employer's face and voice (Mirsky and Lee, 2021; Stupp, 2020). It is also possible to improve the stealth capabilities of attacks by enabling them to proceed without human supervision and aid (making it automatic). For example, if malware could perform a progressive infection of hosts in a network (a.k.a., lateral movement) on its own, then this would reduce command and control (C&C) communication (Truong et al., 2019; Zelinka et al., 2018). Other capabilities include the use of AI to find

* Corresponding author.

E-mail addresses: yisroel@post.bgu.ac.il (Y. Mirsky), ambrademontis@gmail.com (A. Demontis), jaidip@post.bgu.ac.il (J. Kotak), ramk@microsoft.com (R. Shankar), gelei.deng@ntu.edu.sg (D. Gelei), yangliu@ntu.edu.sg (L. Yang), xyzhang@cs.purdue.edu (X. Zhang), maura.pintor@unica.it (M. Pintor), wenke@cc.gatech.edu (W. Lee), elovici@bgu.ac.il (Y. Elovici), battista.biggio@diee.unica.it (B. Biggio).

¹ In this paper, we consider machine learning to be a subset of AI technologies.

zero-day vulnerabilities in software, automate reverse engineering, exploit side channels efficiently, build realistic fake personas, and to perform many more malicious activities with improved efficacy (more examples are presented later in [Section 3](#)).

Offensive Abuse of AI. Adversarial machine learning is the study of security vulnerabilities in AI. It has been shown that an adversary can craft training samples to alter the functionalities of a model e.g., insert a backdoor ([Gu et al., 2017](#)), obtain a desired classification manipulating the test samples (e.g., evade detection) ([Biggio and Roli, 2018](#)) and even infer confidential information about a model ([Orekondy et al., 2019](#)) or the data on which it was trained ([Fredrikson et al., 2015](#)). Since organizations use AI to automate the management, maintenance, operation and defence of their systems and services, an adversary can accomplish their malicious goals by using machine learning *offensively* on these systems (adversarial machine learning).

We note that some attacks are achievable without using or abusing AI. However, attackers can substantially reduce the effort required to perform an attack if they use AI to make it automatic or semi-automatic. By reducing their effort in creating effective strategies, attackers can maximize their return by scaling the attacks in their strength and quantity. Moreover, by acting simultaneously in several phases of the attack chain, the attacker can achieve synergistic effects on the speed and power of the attacks, becoming even more dangerous. On the other hand, some attacks have been enabled by AI, such as the cloning of an individual's voice in a sophisticated social engineering attack ([Brewster, 2021](#)).

1.1. Study overview

In this work, we provide a study of knowledge on offensive AI in the context of enterprise security. The goal of this paper is to help the community (1) better understand the current impact of offensive AI on organizations, (2) prioritize research and development of defensive solutions, and (3) identify trends that may emerge in the near future. This work isn't the first to raise awareness of offensive AI. In [Brundage et al. \(2018\)](#) the authors warned the community that AI can be used for unethical and criminal purposes with examples taken from various domains. In [Caldwell et al. \(2020\)](#) a workshop was held that attempted to identify the potential top threats of AI in criminology. However, both these works relate to the threat of AI on society overall and are not specific to organizations and their networks. Moreover, despite their efforts and preliminary results, these previous analyses provide only examples of how AI can be used to attack and a possible ranking of their risk, while our study gives a structured view of offensive AI through the standard methodologies used to identify potential attack tactics against organizations, deriving strategic insights relevant to defend from these threats.

To accomplish these goals, we performed a literature review to identify the capabilities of an AI-capable adversary. We then performed a panel survey to identify which of these capabilities represent the most relevant threats in practice. There were 35 survey participants: 16 from academia and 19 from industry. The participants from industry were from a wide profile of organizations such as MITRE, IBM, Microsoft, Google, Airbus, Bosch, Fujitsu, Hitachi, and Huawei.

From our literature review, we identified 32 offensive AI capabilities against organizations. Our panel survey revealed that the most significant threats are the capabilities that improve social engineering attacks (e.g., the use of deepfakes to clone the voice of employees). We also found that industry members are most concerned about attacks that enable attackers to steal intellectual property and detect vulnerabilities in their software. Finally, we have also found that modern offensive AI mainly impacts the initial steps of the cyber kill chain (reconnaissance, resource devel-

opment, and initial access). This is because AI technologies are not mature enough to create agents able to carry on attacks that proceed without human supervision and aid. A complete list of our findings can be found in [Section 5.1](#).

1.2. Contributions

In this study, we make the following contributions:

- An overview of how AI can be used to attack organizations and its influence on the cyber kill chain ([Section 2.3](#)).
- An enumeration and description of the 32 offensive AI capabilities that threaten organizations, based on literature review and current events ([Section 3](#)). These capabilities can be categorised as (1) automation, (2) campaign resilience, (3) credential theft, (4) exploit development, (5) information gathering, (6) social engineering, and (7) stealth.
- A threat ranking and insights on how offensive AI impacts organizations, based on a panel survey with members from academia, industry, and government ([Section 4](#)).
- A forecast of the AI threat horizon and the resulting shifts in attack strategies ([Section 5](#)).

1.3. Article structure

This article is structured as follows:

- In [Section 2](#), we provide the reader with a primer on topics which are important for understanding the literature review. The section introduces concepts about AI, offensive AI, and how offensive AI impacts an organization's security.
- In [Section 3](#), we offer our literature review of offensive AI in the context of an organization's security.
- In [Section 4](#), we present the results from a panel survey to help identify the least and most significant threats of offensive AI to organizations.
- In [Section 5](#), we summarize our findings and provide our observations on the matter.

2. Background

In this section, we provide the reader with technical aspects related to offensive AI and introduce offensive AI concepts related to organizations' security. Later in [Section 3](#), we review the latest research on the topic.

2.1. AI And machine learning

AI is a larger domain that mainly deals with creating algorithms that can automate complex tasks. Early AI models were rule-based systems designed using an expert's knowledge ([Yager, 1984](#)), followed by search algorithms for selecting optimal decisions (e.g., finding paths or playing games [Zeng and Church, 2009](#)). Today, the most popular type of AI is machine learning (ML), which is a data-driven approach to AI where programs automatically improve their performance on a task-given experience. Deep learning (DL) is a type of ML where an extensive artificial neural network is used as the predictive model. Breakthroughs in DL have led to its ubiquity in applications such as industrial automation, forecasting, and planning due to its ability to reason upon and generate complex data. Due to the popularity of ML, our literature review inevitably follows this trend. Despite considering all methods and techniques related to using AI in general, we found the vast majority of the offensive AI techniques we found use ML to perform AI-based attacks. Therefore, the majority of the works reviewed in this study involve some form of ML.

Table 1

Examples of where a model can be trained and executed in an attack on an organization. Onsite refers to being within the premises or network of the organization.

Training		Execution		Example
Offsite	Onsite	Offsite	Onsite	
•		•		Vulnerability detection
•			•	Side channel keylogging
	•	•		Channel compression for exfiltration
	•		•	Traffic shaping for evasion
•	•		•	Few-shot learning for record tampering

In general, a machine learning model can be trained on data with explicit ground truth (supervised), with no ground truth (unsupervised), or with a mix of both (semi-supervised). The trade-off between supervised and non-supervised approaches is that supervised methods often have much better performance at a given task but require labeled data which can be expensive or impractical to collect. Moreover, unsupervised techniques are open-world, meaning that they can identify novel patterns that may have been overlooked. Another training paradigm is reinforcement learning, where a model is trained based on reward for good performance. Lastly, for generating content, a popular framework is adversarial learning. This was first popularised in Goodfellow et al. (2014) where the generative adversarial network (GAN) was proposed. A GAN uses a discriminator model to 'help' a generator model produce realistic content by giving feedback on how the content fits a target distribution.

In the context of offensive AI, the location in which an attacker performs training or execution will depend on the attacker's objective and strategy. For example, for reconnaissance tasks, training and execution will likely take place offsite from the organization. However, for attacks, the training and execution may take place onsite, offsite, or both. Another possibility is where the adversary uses few-shot learning (Wang et al., 2020c) by training on general data offsite and then fine tuning on the target data onsite. Additional examples can be found in Table 1. In all cases, the adversary will first design and evaluate their model offsite prior to its usage in the organization to ensure its success and avoid detection.

For onsite execution, an attacker runs the risk of detection if the model is complex (e.g., a DL model). For example, when the model is transferred over to the organization's network or when the attacker's model begins to utilize resources, it may trigger the organization's anomaly detection system. To mitigate this issue, the adversary must consider a trade-off between stealth and effectiveness. For example, the adversary may (1) execute the model during off hours or on non-essential devices, (2) leverage an insider to transfer the model, or (3) transfer the observations off-site for execution.

2.2. Offensive AI

As noted in Section 1, there are two forms of offensive AI (OAI): Attacks using AI and attacks abusing AI. For example, an adversary can (1) use AI to improve the efficiency of an attack (e.g., information gathering, attack automation, and vulnerability discovery) or (2) use knowledge of AI to exploit the defender's AI products and solutions (e.g., to evade a defense or to plant a trojan in a product). The latter form of OAI is commonly referred to as adversarial machine learning.

We will now elaborate on these two forms of offensive AI.

2.2.1. Attacks using AI

Although there are a wide variety of AI tasks which can be used in attacks, we now list the most common ones. Note that these tasks are not mutually exclusive, in fact some build on each

other and produce a synergistic effect on their impact on the attack chain.

Analysis. This is the task of mining or extracting useful insights from data or a model. Some examples of analysis for offense are the use of explainable AI techniques (Ribeiro et al., 2016) to identify how to better hide artifacts (e.g., in malware) and the clustering or embedding of information on an organization to identify assets or targets for social engineering.

Decision Making. The task of producing a strategic plan or coordinating an operation. Examples of this in offensive AI are the use of swarm intelligence to operate an autonomous botnet (Castiglione et al., 2014) and the use of heuristic attack graphs to plan optimal attacks on networks (Bland et al., 2020).

Generation. This is the task of creating content that fits a target distribution which, in some cases, requires realism in the eyes of a human. Examples of generation for offensive uses include the tampering of media evidence (Mirsky et al., 2019; Schreyer et al., 2019), intelligent password guessing (Garg and Ahuja, 2019; Hitaj et al., 2019), and traffic shaping to avoid detection (Han et al., 2020; Novo and Morla, 2020). Deepfakes are another instance of offensive AI in this category. A deepfake is a believable media created by a DL model. The technology can be used to impersonate a victim by puppeting their voice or face to perpetrate a phishing attack (Mirsky and Lee, 2021).

Prediction. This is the task of making a prediction based on previously observed data. Common examples are classification, anomaly detection, and regression. Examples of prediction for an offensive purpose includes the identification of keystrokes on a smartphone based on motion (Hussain et al., 2016; Javed et al., 2020; Marquardt et al., 2011), the selection of the weakest link in the chain to attack (Abid et al., 2018), and the localization of software vulnerabilities for exploitation (Jiang et al., 2019; Lin et al., 2020; Mokhov et al., 2014).

Retrieval. This is the task of finding content that matches or that is semantically similar to a given query. For example, in offense, retrieval algorithms can be used to track an object or an individual in a compromised surveillance system (Rahman et al., 2019; Zhu et al., 2018), to find a disgruntled employee (as a potential insider) using semantic analysis on social media posts, and to summarize lengthy documents (Zhang et al., 2016) during open source intelligence (OSINT) gathering in the reconnaissance phase.

2.2.2. Attacks abusing AI

An attacker can use its AI knowledge to exploit ML model vulnerabilities violating its confidentiality, integrity, or availability (Biggio and Roli, 2018). The vast majority of these attacks is studied in Adversarial Machine Learning, a branch of research that investigates on how to obtain specific malfunctions on ML models to create malicious attacks. These attacks can be staged at either

training (development) or test time (deployment) through one of the following attack vectors:

Modify the Training Data. Here the attacker modifies the training data to harm the integrity or availability of the model. Denial of service (DoS) poisoning attacks (Biggio et al., 2012; Koh and Liang, 2017; Muñoz-González et al., 2017) are when the attacker decreases the model's performance until it is unusable. A backdoor poisoning attack (Chen et al., 2017; Gu et al., 2017) or trojanning attack (Liu et al., 2017), is where the attacker teaches the model to recognize an unusual pattern that triggers a behavior (e.g., classify a sample as safe). A triggerless version of this attack causes the model to misclassify a test sample without adding a trigger pattern to the sample itself (Aghakhani et al., 2021; Shafahi et al., 2018)

Modify the Test Data. In this case, the attacker modifies test samples to have them misclassified (Biggio et al., 2013; Goodfellow et al., 2015; Szegedy et al., 2014). For example, altering the letters of a malicious email to have it misclassified as legitimate, or changing a few pixels in an image to evade facial recognition (Sharif et al., 2016). Therefore, these types of attacks are often referred to as evasion attacks. By modifying test samples ad-hoc to increase the model's resource consumption, the attacker can also slow down the model performances (Shumailov et al., 2021).

Analyze the Model's Responses. Here, the attacker sends a number of crafted queries to the model and observes the responses to infer information about the model's parameters or training data. To learn about the training data, there are membership inference (Shokri et al., 2017), deanonymization (Narayanan and Shmatikov, 2008), and model inversion (Hidano et al., 2017) attacks. For learning about the model's parameters there are model stealing/extraction (Jia et al., 2021; Juuti et al., 2019), and blind-spot detection (Zhang et al., 2019), state prediction (Woh and Lee, 2018).

Modify the Training Code. This is where the attacker performs a supply chain attack by modifying a library used to train ML models (e.g., via an open-source project). For example, compromising a loss (training) function to insert a backdoor (Bagdasaryan and Shmatikov, 2021) or slowing down the created model (Cinà et al., 2022).

Modify the Model's Parameters. In this attack vector, the attacker accesses a trained model (e.g., via a model zoo or security breach) and tampers its parameters to insert a latent behavior. These attacks can be performed at the software (Wang et al., 2020; 2020; Yao et al., 2019) or hardware (Breier et al., 2018a) levels (a.k.a. fault attacks).

Depending on the scenario, an attacker may not have full knowledge or access to the target model:

- **White-Box (Perfect-Knowledge) Attacks:** The attacker knows everything about the target system. This is the worst case for the system defender. Although it is not very likely to happen in practice, this setting is interesting as it provides an empirical upper bound on the attacker's performance.
- **Black-Box (Zero-Knowledge) Attacks:** The attacker knows only the task the model is designed to perform and which kind of features are used by the system in general (e.g., if a malware detector has been trained to perform static or dynamic analysis). The attacker may also be able to analyze the model's responses in a query-based manner to get feedback on certain inputs.
- **Gray-Box (Limited-Knowledge) Attacks:** The attacker has partial knowledge of the target system (e.g., the learning algorithm, architecture, etc.).

In a black or gray box scenario, the attacker can build a surrogate ML model and try to devise the attacks against it as the attacks often transfer between different models. Biggio et al. (2013), Demontis et al. (2019a).

An attacker does not need to be an expert at machine learning to implement these attacks. Many can be acquired from open-source libraries online (Croce and Hein, 2020; Nicolae et al., 2018; Papernot et al., 2018; Pintor et al., 2022).

2.3. Offensive AI vs organizations

In this section, we provide an overview of offensive AI in the context of organizations. First, we review a popular attack model for enterprises. Then we will identify how an AI-capable adversary impacts this model by discussing the adversary's new motivations, goals, capabilities, and requirements. Later in Section 3, we will detail the adversary's techniques based on our literature review.

2.3.1. Attacker motivation

Conventional adversaries use manual effort, common tools, and expert knowledge to reach their goals. In contrast, an AI-capable adversary can use AI to automate its tasks, enhance its tools, and evade detection. These new abilities affect the cyber kill chain.

First, let's discuss why an adversary would consider using AI offensively on an organization. From our literature review (detailed later in Section 3), we observed three reasons why an adversary may be motivated to use offensive AI against an organization: coverage, speed, and success.

Coverage. By using AI, an adversary can scale up its operations by automating complex tasks to decrease human labor and increase the chances of success. For example, AI can be used to automatically craft (Mirsky and Lee, 2021; Stupp, 2020) and launch (employing Leviathan and Matias, 2018; Rebryk and Beliaev, 2020; Singh and Thakur, 2020) spear phishing attacks, distill (Zhang et al., 2016) data collected from OSINT, and reach more assets within a network (Matta et al., 2019; Ou et al., 2005) to gain a stronger foothold. In other words, AI enables adversaries to target more organizations with higher precision attacks with a smaller workforce.

Speed. With AI, an adversary can reach its goals faster. For example, machine learning can be used to help extract credentials (Calzavara et al., 2015; Wang et al., 2019a), intelligently select the next best target during lateral movement (Horák et al., 2019), spy on users to obtain information (e.g., perform speech to text on eavesdropped audio) (White et al., 2011), or find zero-days in software (Jiang et al., 2019; Lin et al., 2020; Mokhov et al., 2014). By reaching a goal faster, the adversary not only saves time for other ventures but can also minimize its presence (duration) within the defender's network.

Success. By enhancing its operations with AI, an adversary increases its likelihood of success. Namely, ML can be used to (1) make the operation more covert by minimizing or camouflaging network traffic (such as C2 traffic) (Han et al., 2020; Novo and Morla, 2020) and by exploiting weaknesses in the defender's AI models such as an ML-based intrusion detection system (IDS) (Sidi et al., 2020), (2) identify opportunities such as good targets for social engineering attacks (Abid et al., 2018) and novel vulnerabilities (Jiang et al., 2019; Lin et al., 2020; Mokhov et al., 2014), (3) enable better attack vectors such as using deepfakes in spear phishing attacks (Stupp, 2020), (4) plan optimal attack strategies (Bland et al., 2020; Horák et al., 2019), and (5) strengthen persistence in the network through automated bot coordination (Castiglione et al., 2014) and malware obfuscation (Datta, 2020).

We note that these motivations are not mutually exclusive. For example, the use of AI to automate a phishing campaign increases coverage, speed, and success.

2.3.2. The attack model

There are a variety of threat agents which target organizations. These agents are cyber terrorists, cyber criminals, employees, hacktivists, nation states, online social hackers, script kiddies, and other organizations like competitors. There are also some non-target specific agents, such as certain botnets and worms, which threaten the security of an organization. A threat agent may be motivated for various reasons. For example, to (1) make money through theft or ransom, (2) gain information through espionage, (3) cause physical or psychological damage for sabotage, terrorism, fame, or revenge, (4) reach another organization, and (5) obtain foothold on the organization as an asset for later use (Krebs, 2014). These agents not only pose a threat to the organization, but also to its employees, customers, and the general public as well (e.g., attacks on critical infrastructure).

In an attack, there may be a number of attack steps that the threat agent must accomplish. These steps depend on the adversary's goal and strategy. For example, in an advanced persistent threat (APT) (Alshamrani et al., 2019; Chen et al., 2018a; Mes-saoud et al., 2016), the adversary may need to reach an asset deep within the defender's network. This would require multiple steps involving reconnaissance, intrusion, lateral movement through a network, and so on. However, some attacks can involve just a single step. For example, a spear phishing attack in which the victim unwittingly provides confidential information or even transfers money. In this paper, we describe the adversary's attack steps using the MITRE ATT&CK Matrix for Enterprise² which captures common adversarial tactics based on real-world observations.

Attacks that involve multiple steps can be thwarted if the defender identifies or blocks the attack early on. The more progress that an adversary makes, the harder it is for the defender to mitigate it. For example, it is better to stop a campaign during the initial intrusion phase than during the lateral movement phase where an unknown number of devices in the network have been compromised. This concept is referred to as the *cyber kill chain*. From an offensive perspective, the adversary will want to shorten and obscure the kill chain to be as efficient and covert as possible. In particular, operation within a defender's network usually requires the attacker to operate through a remote connection or send commands to compromised devices (bots) from a command and control server (C2). This generates presence in the defenders network which can be detected over time.

It is clear that some AI-capable threat agents will be able to perform more sophisticated AI attacks than others. For example, state actors can potentially launch intelligent automated botnets where hacktivists will likely struggle in accomplishing the same. However, we have observed over the years that AI has become increasingly accessible, even to novice users. For example, there are a wide variety of open source deepfakes technologies online which are plug and play³. Therefore, the sophistication gap between certain threat agents may close over time as the availability of AI technology increases.

2.3.3. New threats

AI-capable adversaries have new abilities over conventional cyber adversaries. These abilities give attackers the means to novel acts of sabotage, espionage and theft of intellectual property (IP):

Sabotage. An adversary can use its knowledge to cause damage to an organization in ways that weren't possible before. This is because AI-based adversaries can use (1) adversarial machine learning, (2) generative AI, and (3) deep learning for software analysis.

With adversarial machine learning, an attacker can target the organization's ML products and solutions. For example, they can poison datasets to harm an ML model's performance or plant a backdoor in a model for later exploitation. More examples include, the ability to evade detection in surveillance (Sharif et al., 2016) and affect forecasts models (e.g., finance Goldblum et al., 2021, energy Chen et al., 2019, etc.) With generative AI, an attacker can add or modify evidence in a realistic manner. Examples include the modification of surveillance footage to include or omit evidence (Leetaru, 2019), the tampering of medical scans to harm patients (Mirsky et al., 2019), and the manipulation of financial records to perform fraud (Schreyer et al., 2019). Finally, with recent advances in deep learning, attackers can efficiently and effectively locate vulnerabilities in both source code (Li et al., 2021; 2018) and compiled code (Jiang et al., 2019; Wang et al., 2020b; Xu et al., 2017a). This enables attackers to locate new vulnerabilities for exploitation with minimal effort.

Espionage. With AI, adversaries can spy on organizations in new ways using side-channel analysis and swarm intelligence. Side channels are signals emitted from a device that can be used to infer confidential information (Lavaud et al., 2021). In the past, side-channel attacks were mainly performed in labs using expensive electronics and analytical processes. With AI, adversaries can now perform side-channel attacks on-site and extract information from channels that are temporal, complex, and multi-modal. For example, a compromised smartphone can be used to automatically collect and organize conversations as text using speech-to-text (STT) algorithms, and sentiment analysis (Abd El-Jawad et al., 2018). Attackers can also steal credentials through acoustic and motion side channels (Liu et al., 2015a; Shumailov et al., 2019). AI can also be used to extract latent information from encrypted web traffic (Monaco, 2019), and track users through the organization's social media (Malhotra et al., 2012). Finally, by using swarm intelligence-based malware (Zelinka et al., 2018), attackers can minimize the number of communications that they have to make to maintain and control and progress the attack. Doing so makes it harder for the organization to detect the attacker's presence (i.e., less anomalous outbound traffic) and to remove the malware after blocking the attacker's communication lines.

IP Theft. An AI-capable adversary can extract IP from organizations in new ways. For example, ML models can be stolen from purchased software products, or from cloud services querying the models with crafted inputs (Jia et al., 2021; Juuti et al., 2019). Similar attacks can be performed to steal the model's training data (Haim et al., 2022; Hidano et al., 2017). Obtaining this IP can help an adversary evade or control these models whether they're deployed in the organization or another provider. Another example is AI-based reverse engineering, where compiled software is lifted into higher levels of code so that the algorithms and logic can be understood and stolen (Alrabaei et al., 2021).

2.3.4. OAI attack capabilities

Using the literature review (details later in Section 3), we grouped the papers according to the offensive capability they provide. Doing so revealed 32 offensive AI capabilities (OAC) which di-

² <https://attack.mitre.org/matrices/enterprise/>.

³ <https://github.com/datamillab/awesome-deepfakes-materials>.

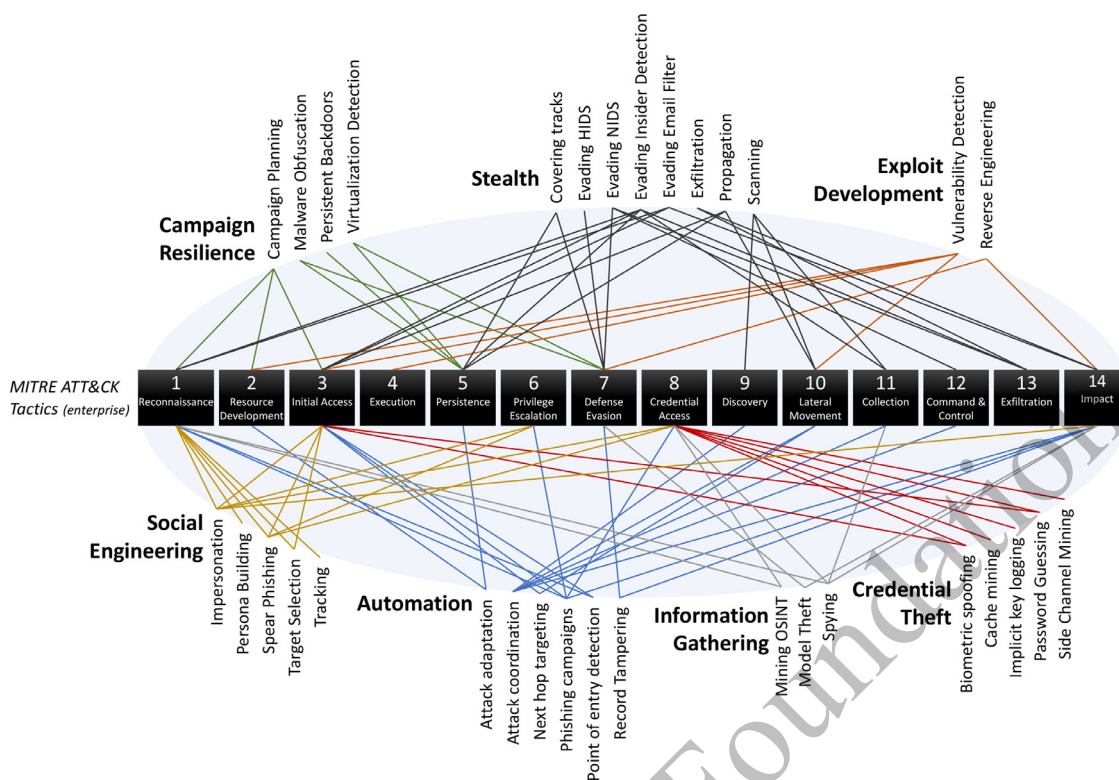


Fig. 1. The 32 offensive AI capabilities (OAC) identified in our literature review, mapped to the MITRE enterprise ATT&CK model. An edge indicates that the OAC directly helps the attacker achieve the indicated attack step.

rectly improve the adversary’s ability to achieve attack steps (e.g., impersonation, user tracking, etc). We then grouped the OACs into categories according to their offensive activity (e.g., social engineering). Finally, we used real use cases reported in the news and by MITRE to validate the OACs and verify that none were missed.

The seven OAC categories were: (1) automation, (2) campaign resilience, (3) credential theft, (4) exploit development, (5) information gathering, (6) social engineering, and (7) stealth. These categories capture the main intent of the adversary reflecting the motivators introduced in Section 2.3.1. Therefore, these categories are non-exclusive (e.g., automating intelligence gathering involves capabilities from both ‘automation’ and ‘information gathering’).

In Fig. 1, we present the OACs and map their influence on the cyber kill chain (the MITRE enterprise ATT&CK model). An edge in the figure means that the indicated OAC improves the attacker’s ability to achieve the given attack step. These edges were obtained by (1) observing real cases reported by MITRE and academic articles and (2) mapping the cases and articles to their respective OACs and their impact on the cyber kill chain. From the figure, we can see that offensive AI impacts every aspect of the attack model. Later in Section 3 we will discuss each of these 32 OACs in greater detail.

These capabilities are materialized in one of two ways:

AI-based tools are programs that perform a specific task in the adversary’s arsenal. For example, a tool for intelligently predicting passwords (Garg and Ahuja, 2019; Hitaj et al., 2019), obfuscating malware code (Datta, 2020), traffic shaping for evasion (Han et al., 2020; Li et al., 2019a; Novo and Morla, 2020), puppeting a persona (Mirsky and Lee, 2021), and so on. These tools are typically in the form of a machine learning model.

AI-driven bots are autonomous bots that can perform one or more attack steps without human intervention, or coordi-

nate with other bots to efficiently reach their goal. These bots may use a combination of swarm intelligence (Castiglione et al., 2014) and machine learning to operate.

3. Literature review

In Section 2.3.4 we presented the 32 offensive AI capabilities. We will now present our literature review of the OACs in order of their 7 categories: automation, campaign resilience, credential theft, exploit development, information gathering, social engineering, and stealth.

Methodology. To perform our literature review, we used the MITRE ATT&CK⁴ matrix as a guide. This matrix lists the common tactics (or attack steps) that an adversary performs when attacking an organization, from planning and reconnaissance leading to the final goal of exploitation. We divided up the work among five different academic workgroups from different international institutions. Each workgroup was assigned a set of tactics from the MITRE ATT&CK matrix, based on their expertise. During the survey, the workgroups were asked to evaluate how AI has been and can be used by an attacker to improve an attacker’s tactics and techniques. Finally, the workgroups cross inspected each other’s content to ensure correctness and completeness.

To identify potential articles and sources to include in our literature review, we selected articles written in the English language and published in peer-reviewed international conference proceedings and journals on the topics of cybersecurity and AI from 1999. As for AI topics, we also included publicly-accessible preprint publications as well since they are well known to be the source of

⁴ <https://attack.mitre.org/matrices/enterprise/>.

the latest advances from key researchers. When searching for attacks which involve AI, we used variations of both 'AI' and 'machine learning' as keywords. The selection process resulted in 225 scientific papers, from which we performed our literature review.

3.1. Automation

The ability to automate complex tasks gives adversaries a hands-off approach to accomplishing attack steps. This not only reduces effort but also increases the adversary's flexibility and enables larger campaigns that are less dependent on C2 signals. Attack automation takes form of either (1) tools which can perform complex tasks using AI (e.g., clone voices, suggest a target) or (2) software (bots) which can operate autonomously to complete an entire attack step with our human intervention (e.g., a bot/malware which propagates on its own by making decisions based on the environment or cooperatively in communication with other bots).

3.1.1. Attack adaptation

Adversaries can use AI to help adapt their malware and attack efforts to unknown environments and find their intended targets. For example, identifying a system (Our, 2020) before attempting an exploit to increase the chances of success and avoid detection. In Black Hat'18, IBM researchers showed how malware can trigger itself using DL by identifying a target's machine by analyzing the victim's face, voice, and other attributes. With models such as decision trees, malware can locate and identify assets via complex rules like (Leong et al., 2019; Lunghi et al., 2017). Instead of transferring screenshots (Arsene, 2020; Brumaghin et al., 2018; Mueller, 2018; Zhang, 2018) DL can be used onsite to extract critical information.

3.1.2. Attack coordination

Cooperative bots can use AI to find the best times and targets to attack. For example, swarm intelligence (Beni, 2020) is the study of autonomous coordination among bots in a decentralized manner. Researchers have proposed that botnets can use swarm intelligence as well. In Zelinka et al. (2018) the authors discuss a hypothetical swarm malware and in Truong et al. (2019) the authors propose another which uses DL to trigger attacks. AI bots can also communicate information on asset locations to fulfill attacks (e.g., send a stolen credential or relevant exploit to a compromised machine).

3.1.3. Next hop targeting

During lateral movement, the adversary must select the next asset to scan or attack. Choosing poorly may prolong the attack and risk detection by the defenders. For example, consider a browser like Firefox which has 4325 key-value pairs denoting the individual configurations. Only some inter-plays of these configurations are vulnerable (Chen et al., 2014; Otsuka et al., 2015). Reinforcement learning can be used to train a detection model which can identify the best browser to target. As for planning multiple steps, a strategy can be formed by using reinforcement learning on Petri nets (Bland et al., 2020) where attackers and defenders are modeled as competing players. Another approach is to use DL (Wu et al., 2021; Yousefi et al., 2018) to explore "attack graphs" (Ou et al. (2005) that contain the target's network structure and the vulnerabilities. Notably, the Q-learning algorithms have enabled the approach to work on large-scale enterprise networks (Matta et al., 2019).

3.1.4. Phishing campaigns

Phishing campaigns involve sending the same emails or robo-phone calls in mass. When someone falls prey and responds, the adversary takes over the conversation. These campaigns can be

fully automated through AI like Google's assistant which can make phone calls on your behalf (Leviathan and Matias, 2018; Rebyrk and Beliaev, 2020; Singh and Thakur, 2020). Furthermore, adversaries can increase their success through mass spear phishing campaigns powered with deepfakes, where (1) a bot calls a colleague of the victim (found via social media), (2) clones his/her voice with 5 seconds of audio (Jia et al., 2018), and then (3) calls the victim in the colleague's voice to exploit their trust.

3.1.5. Point of entry detection

The adversary can use AI to identify and select the best attack vector for an initial infection. For example, in Leslie et al. (2019) statistical models on an organization's attributes were used to predict the number of intrusions it receives. The adversary can train a model on similar information to select the weakest organizations (low-hanging fruits) and the strongest attack vectors.

3.1.6. Record tampering

An adversary may use AI to tamper with records as part of their end goal. For example, ML can be used to impact business decisions with synthetic data (Kumar et al., 2018), to obstruct justice by tampering evidence (Leetaru, 2019), to perform fraud (Schreyer et al., 2019) or to modify medical or satellite imagery (Mirsky et al., 2019). As shown in Mirsky et al. (2019), DL-tampered records can fool human observers and can be accomplished autonomously onsite.

3.2. Campaign resilience

In a campaign, adversaries try to ensure that their infrastructure and tools have a long life. Doing so helps maintain a foothold in the organization and enables the reuse of tools and exploits for future and parallel campaigns. AI can be used to improve campaign resilience through planning, persistence, obfuscation, and detection of virtualization to avoid dynamic analysis.

3.2.1. Campaign planning

Some attacks require careful planning long before the attack campaign to ensure that all of the attacker's tools and resources are obtainable. ML-based cost-benefit analysis tools, such as in Manning et al. (2018), may be used to identify which tools should be developed and how the attack infrastructure should be laid out (e.g., C2 servers, staging areas, etc). It could also be used to help identify other organizations that can be used as beach heads (Krebs, 2014). Moreover, ML can be used to plan a digital twin (Bitton et al., 2018; Fuller et al., 2020) of the victim's network (based on information from reconnaissance) to be created offsite for tuning AI models and developing malware.

3.2.2. Persistent access

An adversary can have bots establish multiple back doors per host and coordinate reinfection efforts among a swarm (Zelinka et al., 2018). Doing so achieves a foothold in an organization by slowing down the effort to purge the campaign. To avoid detection in payloads deployed during boot, the adversary can use a two-step payload that uses ML to identify when to deploy the malware and avoid detection (Anderson, 2017; Fang et al., 2019). Moreover, a USB-sized neural compute stick⁵ can be planted by an insider to enable covert and autonomous onsite DL operations.

⁵ <https://software.intel.com/content/www/us/en/develop/articles/intel-movidius-neural-compute-stick.html>.

3.2.3. Malware obfuscation

ML models such as GANs can be used to obscure a malware's intent from an analyst. Doing so can enable the reuse of the malware, hide the attacker's intents and infrastructure, and prolong an attack campaign. The concept is to take an existing piece of software and emit another piece that is functionally equivalent (similar to translation in NLP). For example, DeepObfusCode (Datta, 2020) uses recurrent neural networks (RNN) to generate ciphered code. Alternatively, backdoors can be planted in open source projects and hidden using similar manners (Pasandi et al., 2019).

3.2.4. Virtualization detection

To avoid dynamic analysis and detection in sandboxes, an adversary may try to have the malware detect the sandbox before triggering. The malware could use ML to detect a virtual environment by measuring system timing (e.g., like in Perianin et al., 2020) and other system properties.

3.3. Credential theft

Although a system may be secure in terms of access control, side channels can be exploited with ML to obtain a user's credentials and vulnerabilities in AI systems can be used to avoid biometric security.

3.3.1. Biometric spoofing

Biometric security is used for access to terminals (such as smartphones) and for performing automated surveillance (Ding et al., 2018; Mozur, 2018; Wang et al., 2017). Recent works have shown how AI can generate "Master Prints" which are deepfakes of fingerprints that can open nearly any partial print scanner (such as on a smartphone) (Bontrager et al., 2018). Face recognition systems can be fooled or evaded with the use of adversarial samples. For example, in Sharif et al. (2016) where the authors generated colorful glasses that alter the perceived identity. Moreover, 'sponge' samples (Shumailov et al., 2021) can be used to slow down a surveillance camera until it is unresponsive or out of batteries (when remote). Voice authentication can also be evaded through adversarial samples, spoofed voice (Wang et al., 2020a), and by cloning the target's voice with deep learning (Wang et al., 2020a).

3.3.2. Cache mining

Information on credentials can be found in a system's cache and log dumps, but a large amount of data makes finding it a difficult task. However, the authors of (Wang et al., 2019a) showed how ML could be used to identify credentials in cache dumps from graphic libraries. Another example is the work of (Calzavara et al., 2015) where an ML system was used to identify cookies containing session information.

3.3.3. Implicit key logging

Over the last few years, researchers have shown how AI can be used as an implicit key-logger by sensing side-channel information from a physical environment. The side channels come in one or a combination of the following aspects:

- Motion.** When tapping on a phone screen or typing on a keyboard, the device and nearby surfaces move and vibrate. Malware can use the smartphone's motion sensors to decipher the touch strokes on the phone (Hussain et al., 2016; Javed et al., 2020) and keystrokes on nearby keyboards (Marquardt et al., 2011). Wearable devices can be exploited in a similar way as well (Liu et al., 2015b; Maiti et al., 2018).
- Audio.** Researchers have shown that, when pressed, each key gives its own unique sound which can be used to infer what is being typed (Compagno et al., 2017; Liu et al., 2015a).

The timing between keystrokes is also a revealing factor due to the structure of the language and keyboard layout. Similar approaches have also been shown for inferring touches on smartphones (Lu et al., 2019; Shumailov et al., 2019; Yu et al., 2019).

Video. In some cases, a nearby smartphone or compromised surveillance camera can be used to observe keystrokes, even when the surface is obscured. For example, via eye movements (Chen et al., 2018b; Wang et al., 2019c; 2018), device motion (Sun et al., 2016), and hand motion (Balagani et al., 2018; Lim et al., 2020).

3.3.4. Password guessing

Humans tend to select passwords with low entropy or with personal information such as dates. GANs can be used to intelligently brute-force passwords by learning from leaked password databases (Hitaj et al., 2019). Researchers have improved on this approach by using RNNs in the generation process (Nam et al., 2020). However, the authors of (Garg and Ahuja, 2019) found that models like (Hitaj et al., 2019) do not work well on Russian passwords. Instead, adversaries may pass the GAN personal information on the user to improve the performance (Seymour and Tully, 2018).

3.3.5. Side channel mining

ML algorithms are adept at extracting latent patterns in noisy data. Adversaries can leverage ML to extract secrets from side channels emitted from cryptographic algorithms. This has been accomplished on a variety of side channels including power consumption (Kocher et al., 1999; Lerman et al., 2014), electromagnetic emanations (Gandolfi et al., 2001), processing time (Brumley and Boneh, 2005), cache hits/misses (Perianin et al., 2020). In general, ML can be used to mine nearly any kind of side channel (Cagli et al., 2017; Heuser et al., 2016; Lerman et al., 2013; Maghrebi et al., 2016; Perin et al., 2021; Picek et al., 2019; 2018; Weissbart et al., 2019). For example, credentials can be extracted from the timing of network traffic (Song et al., 2001).

3.4. Exploit development

Adversaries work hard to understand the content and inner workings of compiled software to (1) steal intellectual property, (2) share trade secrets, (3) and identify vulnerabilities that they can exploit.

3.4.1. Reverse engineering

While interpreting compiled code, an adversary can use ML to help identify functions and behaviors and guide the reversal process. For example, binary code similarity can be used to identify well-known or reused behaviors (Bao et al., 2014; Ding et al., 2019; Duan et al., 2020; Liu et al., 2018; Shin et al., 2015; Xu et al., 2017b; Ye et al., 2020) and autoencoder networks can be used to segment and identify behaviors in code, similar to the work of (Deepreflect, 2021). Furthermore, DL can potentially be used to lift compiled code up to a higher-level representation using graph transformation networks (Yun et al., 2019), similar to semantic analysis in language processing. Protocols and state machines can also be reversed using ML, for example, CAN bus data in vehicles (Huybrechts et al., 2017), network protocols (Li et al., 2015), and commands (Bossert et al., 2014; Wang et al., 2011).

3.4.2. Vulnerability detection

There are a wide variety of software vulnerability detection techniques which can be broken down into static and dynamic approaches:

Static. For open source applications and libraries, the attacker can use ML tools for detecting known types of vulnerabilities in source code (Chakraborty et al., 2020; Feng et al., 2016; Li et al., 2019c; 2018; Mokhov et al., 2014). If its a commercial product (compiled as a binary), then methods such as (Deepreflect, 2021) can be used to identify vulnerabilities by comparing parts of the program's control flow graph to known vulnerabilities.

Dynamic. ML can also be used to perform guided input 'fuzzing' which can reach buggy code faster (Atlidakis et al., 2020; Cheng et al., 2019; Li et al., 2020a; Lin et al., 2020; She et al., 2020; 2019; Wang et al., 2020b). Many works have also shown how AI can mitigate the issue of symbolic execution's massive state space (Janota, 2018; Jiang et al., 2019; Kurin et al., 2019; Liang et al., 2018; Samulowitz and Memisevic, 2007).

3.5. Information gathering

AI scales well and is very good at data mining and language processing. These capabilities can be used by an adversary to collect and distill actionable intel for a campaign.

3.5.1. Mining OSINT

In general, there are three ways in which AI can improve an adversary's OSINT.

Stealth. The adversary can use AI to camouflage its probe traffic to resemble benign services like Google's web crawler (Cohen et al., 2020). Unlike heavy tools like Metagoofil (Martorella, 2020), ML can be used to minimize interactions by prioritizing sites and data elements (Ghazi et al., 2018; Guo et al., 2019).

Gathering. Network structure and elements can be identified using cluster analysis or graph-based anomaly detection (Akoglu et al., 2015). Credentials and asset information can be found using methods like reinforcement learning on other organizations (Schwartz and Kurniawati, 2019). Finally, personnel structure can be extracted from social media using NLP-based web scrappers like Oxylabs (Oxylabs, 2021).

Extracting. Techniques like NLP can be used to translate foreign documents (Dabre et al., 2020), identify relevant documents (Evangelista et al., 2020; Nasar et al., 2019), extract relevant information from online sources (Hlin, 2020; Telegram, 2020), and locate valid identifiers (Malhotra et al., 2012).

3.5.2. Model theft

An adversary may want to steal an AI model to (1) obtain it as intellectual property, (2) extract information about members of its training set (Hidano et al., 2017; Narayanan and Shmatikov, 2008; Shokri et al., 2017), or (3) use it to perform a white-box attack against an organization. As described in Section 2.2.2, if the model can be queried (e.g., model as a service -MAAS), then its parameters (Jia et al., 2021; Juuti et al., 2019) and hyperparameters (Wang and Gong, 2018) can be copied by observing the model's responses. This can also be done through side-channel (Batina et al., 2019) or hardware-level analysis (Breier et al., 2020).

3.5.3. Spying

DL is extremely good at processing audio and video and, therefore, can be used in spyware. For example, a compromised smartphone can map an office by (1) modeling each room with ultrasonic echo responses (Zhou et al., 2017), (2) using object recognition (Jiao et al., 2019) to obtain physical penetration info (control terminals, locks, guards, etc.), and (3) automatically mine relevant information from overheard conversations (Nasar et al., 2019;

Ren et al., 2019). ML can also be used to analyze encrypted traffic. For example it can extract transcripts from encrypted voice calls (White et al., 2011), identify applications (Al-Hababi and Tokgoz, 2020), and reveal internet searches (Monaco, 2019).

3.6. Social engineering

The weakest links in an organization's security are often its humans. Adversaries have long targeted humans by exploiting their emotions and trust. AI provides adversaries with enhanced capabilities to exploit humans further.

3.6.1. Impersonation (Identity theft)

An adversary may want to impersonate someone for a scam, blackmail attempt, defamation attack, or to perform a spear phishing attack with their identity. This can be accomplished using deepfake technologies, which enable the adversary to reenact (puppet) the voice and face of a victim, or alter the existing media content of a victim (Mirsky and Lee, 2021). Recently, the technology has advanced to the state where reenactment can be performed in real-time (Nirkin et al., 2019), and training only requires a few images (Siarohin et al., 2019) or seconds of audio (Jia et al., 2018) from the victim. For high-quality deepfakes, large amounts of audio/video data are still needed. However, when put under pressure, a victim may trust a deepfake even if it has a few abnormalities (e.g., in a phone call) (Workman, 2008). Moreover, the audio/video data may be an end goal inside the organization (e.g., customer data).

3.6.2. Persona building

Adversaries build fake personas on online social networks (OSN) to connect with their targets (Hao, 2019). To evade fake profile detectors, a profile can be cloned and slightly altered using AI (Salminen et al., 2019; 2020; Spiliotopoulos et al., 2020) so that they will appear different yet reflect the same personality. The adversary can then use a number of AI techniques to alter or mask the photos from detection (Li et al., 2019b; Shan et al., 2020; shaoanlu, 2020; Sun et al., 2018). To build connections, a link prediction model can be used to maximize the acceptance rate (Kong and Tong, 2020; Wang et al., 2019b) and a DL chatbot can be used to maintain the conversations (Roller et al., 2021).

3.6.3. Spear phishing

Call-based spear phishing attacks can be enhanced using real-time deepfakes of someone the victim trusts. For example, this occurred in 2019 when a CEO was scammed out of \$240k (Stupp, 2020). For text-based phishing, tweets (zerofox, 2020) and emails (Das and Verma, 2019; Seymour and Tully, 2016; 2018) can be generated to attract a specific victim, or style transfer techniques can be used to mimic a colleague (Fu et al., 2018; Yang et al., 2018).

3.6.4. Target selection

An adversary can use AI to identify victims in the organization who are the most susceptible to social engineering attacks (Abid et al., 2018). It is also possible to build a model based on the target's social attributes (conversations, attended events, etc.) (Bitton et al., 2020; Solomon et al., 2022). Moreover, sentiment analysis can be used to find disgruntled employees to be recruited as insiders (Abd El-Jawad et al., 2018; Dhaoui et al., 2017; Ghiassi and Lee, 2018; Panagiotou et al., 2019; Rathi et al., 2018).

3.6.5. Tracking

To study members of an organization, adversaries may track the member's activities. With ML, an adversary can trace person-

nel across different social media sites by content (Malhotra et al., 2012) and through facial recognition (Black, 2018). ML models can also be used on OSN content to track a member's location (Pellet et al., 2019). Finally, ML can also be used to discover hidden business relationships (Ma et al., 2009; Zhang et al., 2012) from the news and OSNs as well (Kumar and Rathore, 2016; Zhang and Chen, 2018).

3.7. Stealth

In multi-step attacks, covert operations are necessary to ensure success. An adversary can either use or abuse AI to evade detection.

3.7.1. Covering tracks

To hide traces of the adversary's presence, anomaly detection can be performed on the logs to remove abnormal entries (Cao et al., 2017; Debnath et al., 2018). CryptoNets (Gilad-Bachrach et al., 2016) can also be used to hide malware logs and onsite training data for later use. To avoid detection onsite, trojans can be planted in DL intrusion detection systems (IDS) in a supply chain attack at both the hardware (Breier et al., 2018a; 2018b) and software (Li et al., 2022; Liu et al., 2017) levels. DL hardware trojans can use adversarial machine learning to avoid being detected (Hasegawa et al., 2020).

3.7.2. Evading HIDS (Malware detectors)

The struggle between security analysts and malware developers is a never-ending battle, with the malware quickly evolving and defeating detectors. In general, state-of-the-art detectors are vulnerable to evasion (Demontis et al., 2019b; Kolosnjaji et al., 2018; Maiorca et al., 2020). For example, adversaries can evade an ML-based HIDS that performs dynamic analysis by splitting the malware's code into small components executed by different processes (Ispoglou and Payer, 2016). They can also evade ML-based detectors that perform static analysis by adding bytes to the executable (Suciu et al., 2019) or code that does not affect the malware behavior (Anderson et al., 2018; Demetrio et al., 2021; Fang et al., 2019; Pierazzi et al., 2020; Zhiyang et al., 2019). Modifying the malware without breaking its malicious functionality is not easy. Attackers may use AI explanation tools like LIME (Ribeiro et al., 2016) to understand which parts of malware are being recognized by the detector and change them manually. Tools for evading ML-based detection can be found freely online.⁶

3.7.3. Evading NIDS (Network intrusion detection systems)

There are several ways an adversary can use AI to avoid detection while entering, traversing, and communicating over an organization's network. Regarding URL-based NIDSs, attackers can avoid phishing detectors by generating URLs that do not match known examples (Bahnsen et al., 2018). Bots trying to contact their C2 server can generate URLs that appear legitimate to humans (Peck et al., 2019), or that can evade malicious-URL detectors (Sidi et al., 2020). To evade traffic-based NIDSs, adversaries can shape their traffic (Han et al., 2020; Novo and Morla, 2020) or change their timing to hide it (Sharon et al., 2021).

3.7.4. Evading insider detectors

To avoid insider detection mechanisms, adversaries can mask their operations using ML. For example, given one user's credentials, they can use the information on the user's role and the organization's structure to ensure that the operation performed looks legitimate (Sutro, 2020).

3.7.5. Evading email filter

Many email services use machine learning to detect malicious emails. However, adversaries can use adversarial machine learning to evade detection (Dalvi et al., 2004; Gao et al., 2018; Lowd and Meek, 2005a; 2005b). Similarly, malicious documents attached to emails, containing malware, can evade detection as well (e.g., Li et al., 2020b). Finally, an adversary may send emails to be intentionally detected so that they will be added to the defender's training set, as part of a poisoning attack (Biggio et al., 2011).

3.7.6. Exfiltration

Similar to evading NIDSs, adversaries must evade detection when trying to exfiltrate data outside of the network. This can be accomplished by shaping traffic to match the outbound traffic (Li et al., 2019a) or by encoding the traffic within a permissible channel like Facebook chat (Rigaki and Garcia, 2018). To hide the transfer better, an adversary could use DL to compress Patel et al. (2019) and even encrypt (Abadi and Andersen, 2016) the data being exfiltrated. To minimize throughput, audio and video media can be summarized to textual descriptions onsite with ML before exfiltration. Finally, if the network is air-gapped (isolated from the Internet) (Guri and Elovici, 2018) then DL techniques can be used to hide data within side channels such as noise in audio (Jiang et al., 2020).

3.7.7. Propagation & scanning

For stealthy lateral movement, an adversary can configure their Petri nets or attack graphs (see Section 3.1.3) to avoid assets and subnets with certain IDSs and favor networks with more noise to hide in. Moreover, AI can be used to scan hosts and networks covertly by modeling its search patterns and network traffic according to locally observed patterns (Li et al., 2019a).

4. Panel survey & threat ranking

In our literature review (Section 3), we identified the potential offensive AI capabilities (OAC) that an adversary could use to attack an organization. However, some OACs may be impractical, whereas others may pose much larger threats. Therefore, we performed a panel survey to rank these threats and understand their impact on the cyber kill chain.

4.1. Survey setup

We surveyed 35 experts in both subjects of AI and cybersecurity. To be included in the panel survey, a participant must (1) be actively working in academia, industry or government and (1) have at least 2 years experience in both cybersecurity and AI.

From the industry and government sectors, we had 19 participants. Among them were a CISO of a large institution, a CTO and founder of AI-based security companies, an AI ethics researcher from a cybersecurity company, two research managers involved in cyber security AI projects, and seven researchers working in cybersecurity or AI-based cybersecurity. From academia, we had 16 participants: 8 professors and 8 research scientists (Ph.D. and above) with experience in both AI and cyber security. Some of our participants were from MITRE, IBM Research, Microsoft, Airbus, Bosch (RBEI), Fujitsu Ltd., Hitachi Ltd., Huawei Technologies, Nord Security, Institute for Infocomm Research (I2R), Google, Robust Intelligence, Pluribus One, Ermes Cyber Security, Mandiant, WiData, Purdue University, Georgia Institute of Technology, Munich Research Center, University of Cagliari, University of Venice, King's College London, Technische Universität Braunschweig, and the Nanyang Technological University (NTU). The responses of the participants have been anonymized and reflect their own personal views and not the views of their employers.

⁶ https://github.com/zangobot/secml_malware.

The survey consisted of 204 questions that asked the participants to (1) rate different aspects of each OAC, (2) give their opinion on the utility of AI to the adversary in the cyber kill chain, and (3) give their opinion on the balance between the attacker and defender when both have AI. Prior to filling out the questionnaire, all participants were given context of how offensive AI threatens organisations. Prior to rating the aspects of an OAC, participants were given one or more example instances of the OAC for clarification. The questions and the example instances can be found in the appendix. The survey was facilitated using a Google form and it took the participants approximately 30–60 minutes each to complete the form. The responses from the survey were used to produce threat rankings and to gain insights into the threat of offensive AI to organizations.

Only 35 individuals participated in the survey because AI-cybersecurity experts are very busy and hard to reach. However, given the diversity of the participants, we believe that these results still provide meaningful insights into the opinions and concerns that members of academia and industry have on offensive AI.

4.2. Threat ranking

In this section, we measure and rank the various threats of an adversary which can utilize or exploit AI technologies to enhance their attacks. For each OAC the participants were asked to rate four aspects⁷ in the range of 1–7 (low to high):

Profit (P): The amount of benefit that a threat agent gains by using AI compared to using non-AI methods. For example, attack success, flexibility, coverage, attack automation, and persistence. Here profit assumes that the AI tool has already been implemented.

Achievability (A): How easy is it for the attacker to use AI for this task considering that the adversary must implement, train, test, and deploy the AI. This measure also includes the monetary cost to the attacker.

Defeatibility (D): How easy is it for the defender to detect or prevent the AI-based attack. Here, a higher score is bad for the adversary (1=hard to defeat, 7=easy to defeat).

Harm (H): The amount of harm that an AI-capable adversary can inflict in terms of physical, physiological, or monetary damage (including effort put into mitigating the attack).

We say that an adversary is motivated to perform an attack if there is high profit P and high achievability A . Moreover, if there is high P but low A or vice versa, some actors may be tempted to try anyways. Therefore, we model the motivation of using an OAC as $M = \frac{1}{2}(P + A)$. However, just because there is motivation, it does not mean that there is a risk. If the AI attack can be easily detected or prevented, then no amount of motivation will make the OAC a risk. Therefore, we model risk as $R = \frac{M}{D}$ where a low defeatibility (hard to prevent) increases R and a high defeatibility (easy to prevent) lowers R . Risk can also be viewed as the likelihood of the attack occurring, or the likelihood of attack success. Finally, to model threat, we must consider the amount of harm done to the organization. An OAC with high R but no consequences is less of a threat. Therefore, we model our threat score as

$$T = H \frac{\frac{1}{2}(P + A)}{D} = H \frac{M}{D} = HR \quad (1)$$

Before computing T , we normalize P , A , D , and H from the range 1–7 to 0–1. This way, a threat score greater than 1 indicates a significant threat because for these scores (1) the adversary will attempt the attack ($M > D$), and (2) the level of harm will be greater

than the ability to prevent the attack ($\frac{D}{M} < H \leq 1$). We can also see from our model that as an adversary's motivation increases over defeatibility, the amount of harm deemed threatening decreases. This is intuitive because if an attack is easy to achieve and highly profitable, then it will be performed more often. Therefore, even if it is less harmful, attacks will frequently occur so that the damage will be higher in the long run.

4.2.1. OAC threat ranking

In Fig. 2 we present the average P , A , D , and H scores for each OAC. In Fig. 3 we present the OACs ranked according to their threat score T , and contrast their risk scores R to their harm scores H .

The results show that 19 of the OACs (60%) are considered to be significant threats (have a $T > 1$). In general, we observe that the top threats mostly relate to social engineering and malware development. The top three OACs are impersonation, spear phishing, and model theft. These OACs have significantly larger threat scores than the others because they are (1) easy to achieve, (2) have high payoffs, (3) are hard to prevent, and (4) cause the most harm (top left of Fig. 2). Interestingly, the use of AI to run phishing campaigns is considered a large threat even though it has a relatively high D score. We believe this is because, with AI, an adversary can both increase the number and quality of phishing attacks. Therefore, even if 99% of the attempts fail, some will get through and cause the organization damage. The least significant threats were scanning and cache mining which is perceived to have little benefit for the adversary because they pose a high risk of detection. Other low-ranked threats include some on-site automation for propagation, target selection, lateral movement, and covering tracks.

4.2.2. Industry vs academia

In Fig. 4 we look at the average threat scores for each OAC category, and contrast the opinions of members from academia to those from industry.

In general, it seems that academia views AI as a more significant threat to organizations than industry. One can argue that the discrepancy is because industry tends to be more practical and grounded in the present, where academia considers potential threats thus considering the future. For example, when looking at the threat scores from academia, all of the categories are considered significant threats ($T > 1$). However, when looking at the industry's responses, the categories of stealth, credential theft, and campaign resilience are not. This may be because these concepts have presented (proven) themselves less in the wild than the others.

Regardless, both industry and academia seem to agree on the top three most threatening OAC categories: (1) social engineering, (2) information gathering and (3) exploit development. This is because, for these categories, the attacker benefits greatly from using AI (P), can easily implement the relevant AI tools (A), the attack causes considerable damage (H), and there is little the defender can do to prevent them (D) (indicated in Fig. 2). For example, deepfakes are easy to implement yet hard to detect in practice (e.g., in a phone call), and extracting private information from side channels and online resources can be accomplished with little intervention.

Surprisingly, it would appear that both academia and industry consider the use of AI for stealth as the least threatening OAC category in general. Even though there has been a great deal of work showing how IDS models are vulnerable (Novo and Morla, 2020; Suci et al., 2019), IDS evasion approaches were considered the second most defeatable OAC after intelligent scanning. This may have to do with the fact that the adversary cannot evaluate its AI-based evasion techniques inside the actual network and thus risks detection.

Overall, there were some disagreements between our participants from industry and academia regarding the most threatening

⁷ The aspects are based on those proposed by Caldwell et al. (2020).

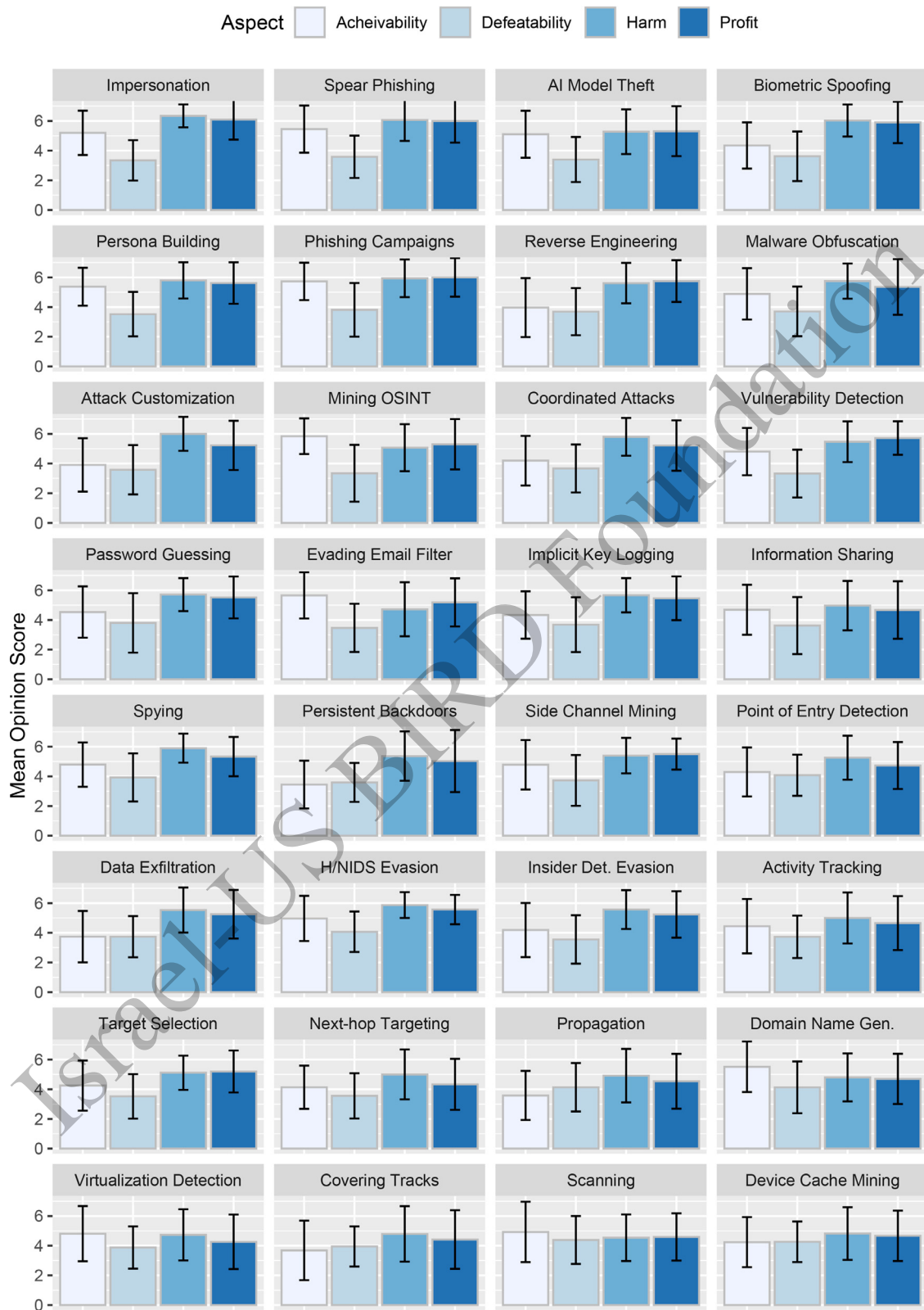


Fig. 2. Survey results: the averaged and normalized opinion scores for each offensive AI capability (OAC) when used against an organization. The OACs are ordered according to their threat score, left to right, starting from the first row.

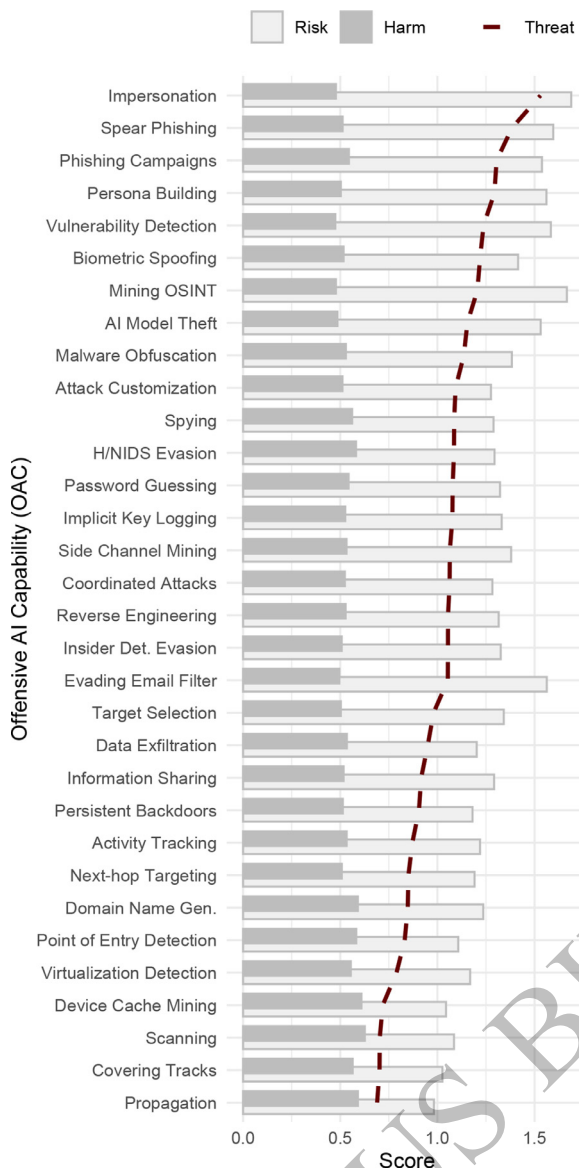


Fig. 3. Survey results: the offensive AI capabilities ranked according to their threat scores.

OACs. The top 10 most threatening OACs for organizations (out of 32) were ranked as follows:

Industry’s Perspective

1. Impersonation
2. Spear Phishing
3. Phishing Campaigns
4. Persona Building
5. Vulnerability Detection
6. Reverse Engineering
7. H/NIDS Evasion
8. Mining OSINT
9. Password Guessing
10. Attack Customization

Academia’s Perspective

1. Impersonation
2. Biometric Spoofing
3. Target Selection
4. Spear Phishing
5. Mining OSINT

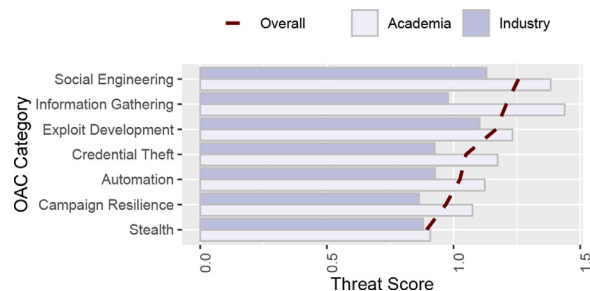


Fig. 4. Survey results: the offensive AI capability categories ranked according to their average threat scores. The scores from industry and academia participants are also presented separately.

6. Vulnerability Detection
7. Spying
8. Persona Building
9. Phishing Campaigns
10. AI Model Theft

Both industry and academia view impersonation as the greatest threat to organizations. This is understandable given recent events where deepfakes were successfully used for impersonation and fraud (FBI, 2022; Fraudster, 2020; Navalny, 2021; Vincent, 2022). We note that our participants from academia view biometric spoofing as the second largest threat, where our participants from industry don’t even consider it in their top 10. We think this is because the latest research on this topic involves ML which can be evaded (e.g., Bontrager et al., 2018; Sharif et al., 2016). In contrast to the academics, our industry participants view this OAC as less harmful to the organization and less profitable to the adversary, perhaps because biometric security is not a common defense used in organization. Regardless, biometric spoofing is still considered the 4-th highest threat overall (Fig. 3). Another insight is that academia is more concerned about the use of ML for spyware, target selection, and the theft of AI models than industry. This may be because these are topics which have long been discussed in academia, but have yet to cause major disruptions in the real-world. For industry, they are more concerned with the use of AI for exploit development, defence evasion and social engineering, likely because these are threats which are out of their control.

4.3. Impact on the cyber kill chain

For each of the 14 MITRE ATT&CK steps, we asked the participants whether they agree or disagree⁸ to the following statements: (1) It more beneficial for the attacker to use AI than conventional methods in this attack step, and (2) AI benefits the attacker more than AI benefits the defender. The objective of these questions were to identify how AI impacts the kill chain and whether AI forms any asymmetry between the attacker and defender.

In Fig. 5 we present the mean opinion scores along with their standard deviations. Overall, our participants felt that AI enhances the adversary’s ability to traverse the kill chain. In particular, we observe that adversary benefits considerably from AI during the first three steps. One explanation is that these attacks are maintained offsite and thus are easier to develop and have less risk. Moreover, we understand from the results that there is a general feeling that defenders do not have a good way to prevent adversarial machine learning attacks. Therefore, AI not only improves defence evasion but also gives the attacker a considerable advantage over the defender in this regard.

⁸ Measured using a 7-step likert scale ranging from strongly disagree (-3) to neutral (0) to strongly agree (+3).

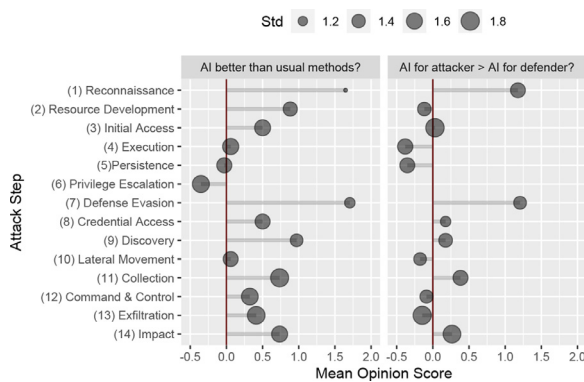


Fig. 5. Survey results: Mean opinion scores on whether (1) it is more beneficial for the adversary to use AI over conventional methods, and (2) AI benefits attackers more than AI benefits defenders. The scores range from -3 to +3.

Our participants also felt that an adversary with AI has a somewhat greater advantage over a defender with AI for most attack steps. In particular, the defender cannot effectively utilize AI to prevent reconnaissance except for mitigating a few kinds of social engineering attacks. Moreover, the adversary has many new uses for AI during the impact step, such as the tampering of records, which the defender does not. However, the participants felt that the defender has an advantage when using AI to detect execution, persistence, and privilege escalation. This is understandable since the defender can train and evaluate models onsite whereas the attacker cannot.

5. Findings & discussion

In this section, we (1) present our main findings from the literature review and panel survey and (2) share our insights on our findings and discuss the road ahead.

5.1. Main findings

From the Literature Review.

- We first observed that there are three primary motivations for an adversary to use AI: coverage, speed, and success (See Section 2.3.1).
- Offensive AI introduces new threats to organizations. A few examples include the poisoning of machine learning models (Biggio et al., 2012; Gu et al., 2017), theft of credentials through side-channel analysis (Song et al., 2001), and the targeting of proprietary training datasets (Hidano et al., 2017; Juuti et al., 2019).
- Adversaries can employ 32 offensive AI capabilities against organizations. These are categorized into seven groups: (1) attack automation, (2) campaign resilience, (3) credential theft, (4) exploit development, (5) information gathering, (6) social engineering, and (7) stealth.
- Defense solutions, such as AI methods for vulnerability detection (Lin et al., 2020), pen-testing (zerofox, 2020), and credential leakage detection (Calzavara et al., 2015) can be weaponized by adversaries for malicious purposes.

From the Panel Survey.

- The top three most threatening categories of offensive AI capabilities against organizations are (1) social engineering, (2) information gathering and (3) exploit development.
- 19 of the 32 offensive AI capabilities pose significant threats to organizations.

- Both industry and academia ranked the threat of using AI for impersonation (e.g., real-time deepfakes to perpetrate phishing and other social engineering attacks) as the highest threat.
- Aside from social engineering aspects, industry and academia are not aligned on the top threats of offensive AI against organizations. Industry members are most concerned with AI being used for reverse engineering, with a focus on the loss of intellectual property and vulnerability detection. Academics, on the other hand, are most concerned about AI being used to perform biometric spoofing (e.g., evading fingerprint and facial recognition) and attack automation.
- Although the evasion of intrusion detection systems (e.g., with adversarial machine learning) is classified as a significant threat, it only ranks number 12 on the list. This may be due to the challenge of the adversary creating effective black box attacks in an unknown IT environment.
- AI impacts the start of the cyber kill chain the most (i.e., reconnaissance, resource development, and initial access). This is because the adversary has more information available and can use this information to refine and evaluate the attacks offsite before proceeding.
- Because AI can be used to automate processes, adversaries may shift from having a few slow covert campaigns to having numerous fast-paced campaigns to overwhelm defenders and increase their chances of success.

5.2. Insights, observations, & limitations

Top Threats. It is understandable why the highest-ranked threats to organizations relate to social engineering attacks and software analysis (vulnerability detection and reverse engineering). It is because these attacks are out of the defender's control. Humans have highly evolved and efficient perception and decision-making abilities. These rely on mental models formed throughout our lives. These mental models (like AI models) can be exploited by presenting information in ways that deceive them (Hollnagel et al., 2006; Woods and Hollnagel, 2006). With deepfakes, social engineering attacks have become even more frequent (Mirsky and Lee, 2021). The same holds for software analysis where ML has been shown to be effective at analyzing software (complex structural data) whether it is source code or a compiled binary (Alrabae et al., 2021; Li et al., 2021; Ye et al., 2020). As mentioned earlier, we believe the reason academia is the most concerned with biometrics is that it almost exclusively uses ML, and academia is well aware of ML's flaws. Industry members may view these attacks as less threatening because physical infiltration is not a top security threat to organizations (Software, 2021). This might explain why they perceive AI attacks on their software and personnel as the greatest threats.

The Near Future. Over the next few years, we believe that there will be an increase in offensive AI incidents, but only at the front and back of the attack model (recon., resource development, and impact – such as record tampering). This is because currently, AI cannot effectively learn new tasks on its own. Therefore, we aren't likely to see botnets that can autonomously and dynamically interact with a diverse set of complex systems (like an organization's network) in the near future. Therefore, since modern adversaries have limited information on the organizations' networks, they are restricted to attacks where the data collection, model development, training, and evaluation occur offsite. In particular, we note that DL models are large and require a considerable amount of resources to run. This makes them easy to detect when transferred into the network or executed onsite. However, the model's footprint might become less anomalous over time as DL proliferates. In the near future, we also expect that phishing campaigns will become more

rampant and dangerous as humans and bots are given the ability to make convincing deepfake phishing calls.

AI is a Double-Edged Sword. We observed that AI technologies for security could also be used in an offensive manner. Some technologies have a dual purpose. For example, ML research into disassembly, vulnerability detection, and penetration testing can be used for both malicious and defensive activities. Some technologies can be repurposed. For example, instead of using explainable AI to validate malware detection, it can be used to hide artifacts (Kuppa and Le-Khac, 2020). And some technologies can be inverted. For example, an insider detection model (Sutro, 2020) can be used to help cover tracks and avoid detection. To help raise awareness, we recommend that researchers note the implications of their work, even for defensive technologies. One caveat is that the usefulness of the ‘sword’ is not symmetric depending on the wielder. For example, generative AI (deepfakes) might be more useful for the attacker because it allows them to generate fake samples (e.g. video) that imitate the benign ones allowing the attacker to accomplish its goal while remaining undetected. Whereas anomaly detection might be more beneficial for the defender.

Limitations of this study. Our study analyzes AI techniques that can be used by attackers against organizations through the MITRE ATT&CK Enterprise matrix. It is also important, however, to note that MITRE also offers other matrices that can be used for different use cases, namely one for Mobile⁹ and one for Industrial Control Systems (ICS).¹⁰ Although the Enterprise and Mobile tactics are almost the same, there are a few unique tactics for ICS that are not contemplated in our study, and that can be extended with the additional non-overlapping threats identified by this scenario.

5.3. The industry's perspective

Using logic to automate attacks is not new to industry – for instance, in 2015, security researchers from FireEye (Intelligence, 2015) found that advanced Russian cyber threat groups built a malware called HAMMERTOSS that used rules based automation to blend its traffic into normal traffic by checking for regular office hours in the time zone and then operating only in that time range. However, the scale and speed that offensive AI capabilities can endow attackers can be damaging.

According to 2019 Verizon Data Breach report analysis of 140 security breaches (Data, 2019), the mean time to compromise an organization and exfiltrating the data ranges is already in the order of minutes. Organizations are already finding automated offensive tactics difficult to combat and anticipate attacks to get stealthier in the future. For instance, according to the final report released by the US National Security Commission on AI in 2021 (Final, 2021a), the warning is clear “The U.S. government is not prepared to defend the United States in the coming artificial intelligence (AI) era.” The final report reasons that this is “Because of AI, adversaries will be able to act with micro-precision, but at macro-scale and with greater speed. They will use AI to enhance cyber attacks and digital disinformation campaigns and to target individuals in new ways.”

Most organizations see offensive AI as an imminent threat – 49% of 102 cybersecurity organizations surveyed by Forrester market research in 2020 (TEOAI, 2020), anticipate offensive AI techniques to manifest in the next 12 months. As a result, more organizations are turning to ways to defend against these attacks. In a 2021 survey (Preparing, 2021b) of 309 organizations' business leaders, C-Suite executives found that 96% of the organizations surveyed are already making investments to guard against AI-powered attacks as they anticipate more automation than what their defenses can handle.

Presently, there are at least three nations which are actively thinking about securing ML systems: The USA through the NSCAI and NIST AI Risk Management, Frameworks¹¹ the UK via their recent release of Principles of securing ML systems,¹² and the EU via the EU AI act coupled with the recently proposed Cyber Resilience Act.¹³ For the most part, these countries emphasise similar aspects: securing the ML pipeline and drawing attention to various attacks on AI systems. It is to be noted that all these frameworks are nascent and are still under discussion. Moreover, their approach is different too. For instance, the NIST framework is voluntary but the proposed EU framework would be mandated for critical ML systems. It is a long road for these standards to come to fruition. Based on followups with our industry members, we believe that organisations may be curious at best about these frameworks but not actively adopting any at this time.

5.4. What's on the horizon

With AI's rapid pace of development and open accessibility, we expect to see a noticeable shift in attack strategies on organizations. First, we foresee that the number of deepfake phishing incidents will increase. In our opinion, this is because the technology (1) is mature, (2) is harder to mitigate than regular phishing, (3) is more effective at exploiting trust, (4) can expedite attacks, and (5) is new as a phishing tactic so cyber defenders are not expecting it. Second, we expect that AI will enable adversaries to target more organizations in parallel and more frequently. As a result, instead of being covert, adversaries may choose to overwhelm the defender's response teams with thousands of attempts for the chance of one success. Finally, as adversaries begin to use AI-enabled bots, defenders will be forced to automate their defenses with bots as well. Keeping humans in the loop to control and determine high-level strategies is a practical and ethical requirement. However, further discussion and research are necessary to form safe and agreeable policies.

5.5. What can be done?

Attacks Using AI. Industry and academia should focus on developing solutions for mitigating the top threats. Personnel can be shown what to expect from AI-powered social engineering and further research can be done on detecting deepfakes, but in a manner that is robust to a dynamic adversary (Mirsky and Lee, 2021). Moreover, we recommend research into post-processing tools that can protect software from analysis after development (i.e., anti-vulnerability detection).

Attacks Against AI. The advantages and vulnerabilities of AI have profoundly questioned their widespread adoption, especially in mission-critical and cybersecurity-related tasks. In the meantime, organizations are working on automating the development and operations of ML models (MLOps), without focusing too much on ML security-related issues. To bridge this gap, we argue that extending the current MLOps paradigm to also encompass ML security (MLSecOps) may be a relevant way toward improving the security posture of such organizations. To this end, we envision the incorporation of security testing, protection and monitoring of AI/ML models into MLOps. Doing so will enable organizations to seamlessly deploy and maintain more secure and reliable AI/ML models.

⁹ <https://attack.mitre.org/techniques/mobile/>.

¹⁰ <https://attack.mitre.org/techniques/ics/>.

¹¹ <https://www.nist.gov/itl/ai-risk-management-framework>.

¹² <https://www.ncsc.gov.uk/collection/machine-learning>.

¹³ <https://digital-strategy.ec.europa.eu/en/library/cyber-resilience-act>.

6. Conclusion

In this study we first explored, categorized, and identified the threats of offensive AI against organizations (Sections 2 and 2.3). We then detailed the threats and ranked them through a panel survey with experts from the domain (Sections 3 and 4). Finally, we provided insights into our results and gave directions for future work (Section 5). We hope this study will be meaningful and helpful to the community in addressing the imminent threat of offensive AI.

Declaration of Competing Interest

None.

Data Availability

No data was used for the research described in the article.

Acknowledgments

The authors would like to thank Laurynas Adomaitis, Sin G. Teo, Manojkumar Parmar, Charles Hart, Matilda Rhode, Dr. Daniele Sgandurra, Dr. Pin-Yu Chen, Evan Downing, Didier Contis, Marco Uras, Konrad Rieck, Dr. Armin Wasicek, Dr. Stefano Traverso, and Dr. Fabio Pierazzi for taking the time to participate in our panel survey. We note that the views reflect the participant's personal experiences and does not reflect the view of the participant's employer. This material is based upon work supported by the Zuckerman STEM Leadership Program. This work has also been partially supported also by the PRIN 2017 project RexLearn (grant no. 2017TWNMH2), funded by the Italian Ministry of Education, University and Research; by the U.S.-Israel Energy Center managed by the Israel-U.S. Binational Industrial Research and Development (BIRD) Foundation; and by BMK, BMDW, and the Province of Upper Austria in the frame of the COMET Programme managed by FFG in the COMET Module S3AI.

Appendix A. The complete questionnaire

A1. Rating the threat

In an attack on an organization, there are 7 malicious activities that can be enhanced using AI: automation, information gathering, campaign resilience, credential theft, social engineering, stealth, and exploit development.

Please rate accordingly:

Harm. How harmful is an attacker with AI in this task? (damage, attack persistence, evasion, defense effort)

Profit. How beneficial is AI to the attacker in this task? (compared to using non-AI methods) (attack success, flexibility, coverage, automation, and persistence). Assume that the AI tool has already been implemented.

Achievability How easy is it for the attacker to use AI for this task? (implement, train, and deploy the AI)

Defeatibility How easy is it for the defender to detect or prevent it? (1=hard to defeat, 7=easy to defeat)

Activity: Automation.

Attack Customization (e.g., adjusting an exploit) (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatibility: __

Coordinated Attacks: (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatibility: __

Information Sharing (among bots or threat agents) (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatibility: __

Next-hop Targeting (e.g., lateral movement) (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatibility: __

Phishing Campaigns (e.g., automated into collection crafting of spear phishing emails, calls, ...)

Harm: __, Profit: __, Achievability: __, Defeatibility: __

Point of Entry Detection (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatibility: __

Activity: Information Gathering (IG).

Mining OSINT (e.g., parsing websites, retrieving relevant info, ...) (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatibility: __

AI Model Theft (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatibility: __

Spying (e.g., collecting and mining conversations from the microphone, locations from the camera,...)

Harm: __, Profit: __, Achievability: __, Defeatibility: __

Activity: Campaign Resilience (CR).

Malware Obfuscation (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatibility: __

Persistent Backdoors (e.g., automated reinfection, backdoor info shared among bots, ...)

Harm: __, Profit: __, Achievability: __, Defeatibility: __

Virtualization Detection (anti-forensics for malware) (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatibility: __

Activity: Credential Theft (CT).

Biometric Spoofing (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatibility: __

Device Cache Mining (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatibility: __

Implicit Key Logging (e.g., using smartphone acceleration, keystroke sounds, ...)

Harm: __, Profit: __, Achievability: __, Defeatibility: __

Intelligent Password Guessing (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatibility: __

Side Channel Mining (e.g., memory or timing patterns) (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatibility: __

Activity: Social Engineering (SE).

Impersonation (e.g., voice, text, video deepfakes and online social profiles) (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatibility: __

Persona Building (e.g., a targeted trustworthy/attractive online profile) (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatibility: __

Spear Phishing (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatability: __

Target Selection (e.g., weakest link with asset) (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatability: __

Activity Tracking (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatability: __

Activity: Stealth.

Covering Tracks (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatability: __

Web Domain Name Generation (e.g., DGAs to avoid detection and blacklisting)

Harm: __, Profit: __, Achievability: __, Defeatability: __

Evading Network or Host-based Intrusion Detection Systems (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatability: __

Evading Insider Detection Systems (e.g., replicate access pattern of other user)

Harm: __, Profit: __, Achievability: __, Defeatability: __

Evading Email Filter (i.e., for SPAM and phishing) (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatability: __

Data Exfiltration (e.g., evading firewall or over an air-gap for an isolated network)

Harm: __, Profit: __, Achievability: __, Defeatability: __

Propagation (lateral movement over a network) (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatability: __

Scanning (e.g., local host, network assets, ports, vulnerabilities, ...) (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatability: __

Activity: Exploit Development (ED).

Reverse Engineering (i.e., to assist in manually finding a vulnerability or steal IP) (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatability: __

Vulnerability Detection (e.g., intelligent fuzzing, static analysis, ...) (1 = low, 7 = high)

Harm: __, Profit: __, Achievability: __, Defeatability: __

A2. The impact on the cyber kill chain

In an advanced persistent threat (APT) an adversary follows 14 tactics to attack an organization according to the MITRE A&TACK matrix. However, at each step the defender can stop the attack and effectively kill the chain of events, preventing the attacker from reaching its goal.

Compared to using conventional methods, AI helps the attacker in...

(strongly disagree, disagree somewhat disagree, neutral, somewhat agree, agree, strongly agree)

(1) Reconnaissance: __, (2) Resource Development: __, (3) Initial Access: __, (4) Execution: __, (5) Persistence: __, (6) Privilege Escalation: __, (7) Defense Evasion: __, (8) Credential Access: __, (9)

Discovery: __, (10) Lateral Movement: __, (11) Collection: __, (12) Command & Control: __, (13) Exfiltration: __, (14) Impact: __

For each tactic, would AI help the attacker more than the defender?

(strongly disagree, disagree somewhat disagree, neutral, somewhat agree, agree, strongly agree)

(1) Reconnaissance: __, (2) Resource Development: __, (3) Initial Access: __, (4) Execution: __, (5) Persistence: __, (6) Privilege Escalation: __, (7) Defense Evasion: __, (8) Credential Access: __, (9) Discovery: __, (10) Lateral Movement: __, (11) Collection: __, (12) Command & Control: __, (13) Exfiltration: __, (14) Impact: __

References

- Abadi, M., Andersen, D. G., 2016. Learning to protect communications with adversarial neural cryptography. [arXiv:1610.06918](https://arxiv.org/abs/1610.06918).
- Abd El-Jawad, M.H., Hodhod, R., Omar, Y.M.K., 2018. Sentiment analysis of social media networks using machine learning. In: 2018 14th International Computer Engineering Conference (ICENCO), pp. 174–176. doi:10.1109/ICENCO.2018.8636124.
- Abid, Y., Imine, A., Rusinowitch, M., 2018. Sensitive attribute prediction for social networks users. *EDBT/ICDT Workshops*.
- Aghakhani, H., Meng, D., Wang, Y., Kruegel, C., Vigna, G., 2021. Bullseye polytope: a scalable clean-label poisoning attack with improved transferability. In: IEEE European Symposium on Security and Privacy, EuroS&P 2021, Vienna, Austria, September 6–10, 2021. IEEE, pp. 159–178. doi:10.1109/EuroSP51992.2021.00021.
- Akoglu, L., Tong, H., Koutra, D., 2015. Graph based anomaly detection and description: a survey. *Data Min. Knowl. Discov.* 29 (3), 626–688.
- Al-Hababi, A., Tokgoz, S.C., 2020. Man-in-the-middle attacks to detect and identify services in encrypted network flows using machine learning. In: 2020 3rd International Conference on Advanced Communication Technologies and Networking (CommNet). IEEE, pp. 1–5.
- Alrabaei, S., Choo, K.-K.R., Qbea'h, M., Khasawneh, M., 2021. BinDeep: binary to source code matching using deep learning. In: 2021 IEEE 20th International Conference on Trust, Security and Privacy in Computing and Communications (TrustCom). IEEE, pp. 1100–1107.
- Alshamrani, A., Myneni, S., Chowdhary, A., Huang, D., 2019. A survey on advanced persistent threats: techniques, solutions, challenges, and research opportunities. *IEEE Commun. Surv. Tutor.* 21 (2), 1851–1877.
- Anderson, H., 2017. Evading machine learning malware detection.
- Anderson, H. S., Kharkar, A., Filar, B., Evans, D., Roth, P., 2018. Learning to evade static pe machine learning malware models via reinforcement learning. 1801.08917.
- Arsene, L., 2020. Oil & gas spearphishing campaigns drop agent tesla spyware in advance of historic opec+ deal. <https://labs.bitdefender.com/2020/04/oil-gas-spearphishing-campaigns-drop-agent-tesla-spyware-in-advance-of-historic-opec-deal/>.
- Atlidakis, V., Geambasu, R., Godefroid, P., Polishchuk, M., Ray, B., 2020. Pythia: grammar-based fuzzing of rest APIs with coverage-guided feedback and learning-based mutations. [arXiv preprint arXiv:2005.11498](https://arxiv.org/abs/2005.11498).
- Bagdasaryan, E., Shmatikov, V., 2021. Blind backdoors in deep learning models. In: 30th USENIX Security Symposium (USENIX Security 21), pp. 1505–1521.
- Bahnsen, A.C., Torroledo, I., Camacho, L.D., Villegas, S., 2018. DeepPhish: Simulating malicious AI. In: 2018 APWG Symposium on Electronic Crime Research (eCrime), pp. 1–8.
- Balagani, K.S., Conti, M., Gasti, P., Georgiev, M., Gurtler, T., Lain, D., Miller, C., Moilas, K., Samarin, N., Saraci, E., et al., 2018. SILK-TV: secret information leakage from keystroke timing videos. In: European Symposium on Research in Computer Security. Springer, pp. 263–280.
- Bao, T., Burket, J., Woo, M., Turner, R., Brumley, D., 2014. {BYTEWEIGHT}: learning to recognize functions in binary code. In: 23rd {USENIX} Security Symposium ({USENIX} Security 14), pp. 845–860.
- Batina, L., Bhasin, S., Jap, D., Picek, S., 2019. CSI NN: reverse engineering of neural network architectures through electromagnetic side channel. In: 28th USENIX Security Symposium (USENIX Security 19). USENIX Association, Santa Clara, CA, pp. 515–532. <https://www.usenix.org/conference/usenixsecurity19/presentation/batina>
- Beni, G., 2020. Swarm intelligence. In: Complex Social and Behavioral Systems: Game Theory and Agent-Based Models, pp. 791–818.
- Biggio, B., Corona, I., Fumera, G., Giacinto, G., Roli, F., 2011. Bagging classifiers for fighting poisoning attacks in adversarial classification tasks. In: Sansone, C., Kittler, J., Roli, F. (Eds.), 10th International Workshop on Multiple Classifier Systems (MCS). Springer-Verlag, pp. 350–359.
- Biggio, B., Corona, I., Maiorca, D., Nelson, B., Šrđić, N., Laskov, P., Giacinto, G., Roli, F., 2013. Evasion attacks against machine learning at test time. In: Blockeel, H., Kersting, K., Nijssen, S., Železný, F. (Eds.), Machine Learning and Knowledge Discovery in Databases (ECML PKDD), Part III. Springer Berlin Heidelberg, pp. 387–402.
- Biggio, B., Nelson, B., Laskov, P., 2012. Poisoning attacks against support vector machines. In: Langford, J., Pineau, J. (Eds.), 29th Int'l Conf. on Machine Learning. Omnipress, pp. 1807–1814.
- Biggio, B., Roli, F., 2018. Wild patterns: ten years after the rise of adversarial machine learning. *Pattern Recognit.* 84, 317–331.

- Bitton, R., Boymgold, K., Puzis, R., Shabtai, A., 2020. Evaluating the information security awareness of smartphone users. In: Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems, pp. 1–13.
- Bitton, R., Gluck, T., Stan, O., Inokuchi, M., Ohta, Y., Yamada, Y., Yagyu, T., Elovici, Y., Shabtai, A., 2018. Deriving a cost-effective digital twin of an ICS to facilitate security evaluation. In: European Symposium on Research in Computer Security. Springer, pp. 533–554.
- Bland, J.A., Petty, M.D., Whitaker, T.S., Maxwell, K.P., Cantrell, W.A., 2020. Machine learning cyberattack and defense strategies. *Comput. Secur.* 92, 101738. doi:10.1016/j.cose.2020.101738.
- Black Hat USA, 2018. <https://www.blackhat.com/us-18/arsenal.html#social-mapper-social-media-correlation-through-facial-recognition>.
- Bontrager, P., Roy, A., Togelius, J., Memon, N., Ross, A., 2018. DeepMasterPrints: generating masterprints for dictionary attacks via latent variable evolution. In: 2018 IEEE 9th International Conference on Biometrics Theory, Applications and Systems (BTAS). IEEE, pp. 1–9.
- Bossert, G., Guihéry, F., Hiet, G., 2014. Towards automated protocol reverse engineering using semantic information. In: Proceedings of the 9th ACM Symposium on Information, Computer and Communications Security, pp. 51–62.
- Breier, J., Hou, X., Jap, D., Ma, L., Bhasin, S., Liu, Y., 2018. Practical fault attack on deep neural networks. In: Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security, pp. 2204–2206.
- Breier, J., Hou, X., Jap, D., Ma, L., Bhasin, S., Liu, Y., 2018. Practical fault attack on deep neural networks. In: Lie, D., Mannan, M., Backes, M., Wang, X. (Eds.), Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security, CCS 2018, Toronto, ON, Canada, October 15–19, 2018. ACM, pp. 2204–2206. doi:10.1145/3243734.3278519.
- Breier, J., Jap, D., Hou, X., Bhasin, S., Liu, Y., 2020. SNIFF: reverse engineering of neural networks with fault attacks. *arXiv preprint arXiv:2002.11021*.
- Brewster, T., 2021. Fraudsters cloned company director's voice in \$35 million bank heist, police find. <https://www.forbes.com/sites/thomasbrewster/2021/10/14/huge-bank-fraud-uses-deep-fake-voice-tech-to-steal-millions/?sh=2e3ea2297559>, (Accessed on 06/15/2022).
- Brumaghin, E., Unterbrink, H., Tacheau, E., 2018. Old dog, new tricks - analysing new RTF-based campaign distributing Agent Tesla, Loki with PyREbox. https://blog.talosintelligence.com/2018/10/old-dog-new-tricks-analysing-new-rtf_15.html.
- Brumley, D., Boneh, D., 2005. Remote timing attacks are practical. *Comput. Netw.* 48 (5), 701–716.
- Brundage, M., Avin, S., Clark, J., Toner, H., Eckersley, P., Garfinkel, B., Dafoe, A., Scharre, P., Zeitoff, T., Filar, B., et al., 2018. The malicious use of artificial intelligence: forecasting, prevention, and mitigation. *arXiv preprint arXiv:1802.07228*.
- Cagli, E., Dumas, C., Prouff, E., 2017. Convolutional neural networks with data augmentation against jitter-based countermeasures. In: International Conference on Cryptographic Hardware and Embedded Systems. Springer, pp. 45–68.
- Caldwell, M., Andrews, J., Tanay, T., Griffin, L., 2020. Ai-enabled future crime. *Crim. Sci.* 9 (1), 1–13.
- Calzavara, S., Tolomei, G., Casini, A., Bugliesi, M., Orlando, S., 2015. A supervised learning approach to protect client authentication on the web. *ACM Trans. Web* 9 (3). doi:10.1145/2754933.
- Cao, Q., Qiao, Y., Lyu, Z., 2017. Machine learning to detect anomalies in web log analysis. In: 2017 3rd IEEE International Conference on Computer and Communications (ICCC), pp. 519–523.
- Castiglione, A., Prisco, R.D., Santis, A.D., Fiore, U., Palmieri, F., 2014. A botnet-based command and control approach relying on swarm intelligence. *J. Netw. Comput. Appl.* 38, 22–33. doi:10.1016/j.jnca.2013.05.002. <https://www.sciencedirect.com/science/article/pii/S1084804513001161>
- Chakraborty, S., Krishna, R., Ding, Y., Ray, B., 2020. Deep learning based vulnerability detection: are we there yet? *arXiv preprint arXiv:2009.07235*.
- Chen, J., Su, C., Yeh, K.-H., Yung, M., 2018a. Special issue on advanced persistent threat.
- Chen, W., Qiao, X., Wei, J., Zhong, H., Huang, X., 2014. Detecting inter-component configuration errors in proactive: a relation-aware method. In: 2014 14th International Conference on Quality Software. IEEE, pp. 184–189.
- Chen, X., Liu, C., Li, B., Lu, K., Song, D., 2017. Targeted backdoor attacks on deep learning systems using data poisoning. *arXiv preprint arXiv:1712.05526*.
- Chen, Y., Li, T., Zhang, R., Zhang, Y., Hedgpeth, T., 2018. EyeTell: video-assisted touchscreen keystroke inference from eye movements. In: 2018 IEEE Symposium on Security and Privacy (SP). IEEE, pp. 144–160.
- Chen, Y., Tan, Y., Zhang, B., 2019. Exploiting vulnerabilities of load forecasting through adversarial attacks. In: Proceedings of the Tenth ACM International Conference on Future Energy Systems, pp. 1–11.
- Cheng, L., Zhang, Y., Zhang, Y., Wu, C., Li, Z., Fu, Y., Li, H., 2019. Optimizing seed inputs in fuzzing with machine learning. In: 2019 IEEE/ACM 41st International Conference on Software Engineering: Companion Proceedings (ICSE-Companion). IEEE, pp. 244–245.
- Cinà, A. E., Demontis, A., Biggio, B., Roli, F., Pelillo, M., 2022. Energy-latency attacks via sponge poisoning. *arXiv:2203.08147 [cs]*.
- Cohen, D., Mirsky, Y., Kamp, M., Martin, T., Elovici, Y., Puzis, R., Shabtai, A., 2020. DANTE: a framework for mining and monitoring darknet traffic. In: European Symposium on Research in Computer Security. Springer, pp. 88–109.
- Compagn, A., Conti, M., Lain, D., Tsudik, G., 2017. Don't skype & type! Acoustic eavesdropping in voice-over-ip. In: Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security, pp. 703–715.
- Croce, F., Hein, M., 2020. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. *ICML*.
- Dabre, R., Chu, C., Kunchukuttan, A., 2020. A survey of multilingual neural machine translation. *ACM Comput. Surv. (CSUR)* 53 (5), 1–38.
- Dalvi, N., Domingos, P., Mausam, Sanghai, S., Verma, D., 2004. Adversarial classification. In: Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD). Seattle, pp. 99–108.
- Das, A., Verma, R., 2019. Automated email generation for targeted attacks using natural language. 1908.06893.
- Datta, S., 2020. DeepObfuscate: source code obfuscation through sequence-to-sequence networks. *arXiv preprint arXiv:1909.01837*.
- , 2019. 2019 Data Breach Investigations Report. Verizon, Inc.
- Debnath, B., Solaimani, M., Gulzar, M.A.G., Arora, N., Lumezanu, C., Xu, J., Zong, B., Zhang, H., Jiang, G., Khan, L., 2018. LogLens: a real-time log analysis system. In: 2018 IEEE 38th International Conference on Distributed Computing Systems (ICDCS), pp. 1052–1062.
- , 2021. DeepReflect: discovering malicious functionality through binary reconstruction. 30th USENIX Security Symposium (USENIX Security 21). USENIX Association. <https://www.usenix.org/conference/usenixsecurity21/presentation/downing>
- Demetrio, L., Biggio, B., Lagorio, G., Roli, F., Armando, A., 2021. Functionality-preserving black-box optimization of adversarial windows malware. *IEEE Trans. Inf. Forensics Secur.* 16, 3469–3478. doi:10.1109/TIFS.2021.3082330.
- Demontis, A., Melis, M., Pintor, M., Jagielski, M., Biggio, B., Oprea, A., Nita-Rotaru, C., Roli, F., 2019. Why do adversarial attacks transfer? Explaining transferability of evasion and poisoning attacks. 28th USENIX Security Symposium (USENIX Security 19). USENIX Association.
- Demontis, A., Melis, M., Pintor, M., Jagielski, M., Biggio, B., Oprea, A., Nita-Rotaru, C., Roli, F., 2019. Why do adversarial attacks transfer? Explaining transferability of evasion and poisoning attacks. 28th USENIX Security Symposium (USENIX Security 19). USENIX Association.
- Dhaoui, C., Webster, C.M., Tan, L.P., 2017. Social media sentiment analysis: lexicon versus machine learning. *J. Consum. Mark.*
- Ding, C., Huang, K., Patel, V.M., Lovell, B.C., 2018. Special issue on video surveillance-oriented biometrics. *Pattern Recognit. Lett.* 107, 1–2.
- Ding, S.H., Fung, B.C., Charland, P., 2019. Asm2Vec: boosting static representation robustness for binary clone search against code obfuscation and compiler optimization. In: 2019 IEEE Symposium on Security and Privacy (SP). IEEE, pp. 472–489.
- Duan, Y., Li, X., Wang, J., Yin, H., 2020. DeepBinDiff: learning program-wide code representations for binary diffing. In: Proceedings of the 27th Annual Network and Distributed System Security Symposium (NDSS'20).
- Evangelista, J.R.G., Sassi, R.J., Romero, M., Napolitano, D., 2020. Systematic literature review to investigate the application of open source intelligence (OSINT) with artificial intelligence. *J. Appl. Secur. Res.* 1–25.
- Fang, Z., Wang, J., Li, B., Wu, S., Zhou, Y., Huang, H., 2019. Evading anti-malware engines with deep reinforcement learning. *IEEE Access* 7, 48867–48879. doi:10.1109/ACCESS.2019.2908033.
- FBI, 2022. FBI: Scammers are interviewing for remote jobs using deepfake tech - mashable. <https://mashable.com/article/deepfake-job-interviews-fbi#:~:text=Deepfakes%20involve%20using%20AI%2Dpowered,say%20whatever%20you'd%20like>, (Accessed on 08/17/2022).
- Feng, Q., Zhou, R., Xu, C., Cheng, Y., Testa, B., Yin, H., 2016. Scalable graph-based bug search for firmware images. In: Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, pp. 480–491.
- , 2021. Final Report. National Security Commission on Artificial Intelligence.
- Fraudsters cloned company director's voice in \$35 million bank heist, police find. <https://www.forbes.com/sites/thomasbrewster/2021/10/14/huge-bank-fraud-uses-deep-fake-voice-tech-to-steal-millions/?sh=23254d367559>, (Accessed on 08/17/2022).
- Fredrikson, M., Jha, S., Ristenpart, T., 2015. Model inversion attacks that exploit confidence information and basic countermeasures. In: Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security. ACM, New York, NY, USA, pp. 1322–1333.
- Fu, Z., Tan, X., Peng, N., Zhao, D., Yan, R., 2018. Style transfer in text: exploration and evaluation. In: McIlraith, S.A., Weinberger, K.Q. (Eds.), Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2–7, 2018. AAAI Press, pp. 663–670. <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/17015>
- Fuller, A., Fan, Z., Day, C., Barlow, C., 2020. Digital twin: enabling technologies, challenges and open research. *IEEE Access* 8, 108952–108971.
- Gandolfi, K., Mourtel, C., Olivier, F., 2001. Electromagnetic analysis: concrete results. In: International Workshop on Cryptographic Hardware and Embedded Systems. Springer, pp. 251–261.
- Gao, J., Lanchantin, J., Soffa, M.L., Qi, Y., 2018. Black-box generation of adversarial text sequences to evade deep learning classifiers. In: 2018 IEEE Security and Privacy Workshops (SPW), pp. 50–56. doi:10.1109/SPW.2018.00016.
- Garg, V., Ahuja, L., 2019. Password guessing using deep learning. In: 2019 2nd International Conference on Power Energy, Environment and Intelligent Control (PEEIC). IEEE, pp. 38–40.
- Ghazi, Y., Anwar, Z., Mumtaz, R., Saleem, S., Tahir, A., 2018. A supervised machine learning based approach for automatically extracting high-level threat intelligence from unstructured sources. In: 2018 International Conference on Frontiers of Information Technology (FIT). IEEE, pp. 129–134.

- Ghiassi, M., Lee, S., 2018. A domain transferable lexicon set for twitter sentiment analysis using a supervised machine learning approach. *Expert Syst. Appl.* 106, 197–216.
- Gilad-Bachrach, R., Dowlin, N., Laine, K., Lauter, K., Naehrig, M., Wernsing, J., 2016. CryptoNets: applying neural networks to encrypted data with high throughput and accuracy. In: *International Conference on Machine Learning*, pp. 201–210.
- Goldblum, M., Schwarzschild, A., Patel, A., Goldstein, T., 2021. Adversarial attacks on machine learning systems for high-frequency trading. In: *Proceedings of the Second ACM International Conference on AI in Finance*, pp. 1–9.
- Goodfellow, I.J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y., 2014. Generative adversarial nets. In: *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2*. MIT Press, Cambridge, MA, USA, pp. 2672–2680.
- Goodfellow, I.J., Shlens, J., Szegedy, C., 2015. Explaining and harnessing adversarial examples. In: *International Conference on Learning Representations*.
- Gu, T., Dolan-Gavitt, B., Garg, S., 2017. BadNets: identifying vulnerabilities in the machine learning model supply chain. *NIPS Workshop on Machine Learning and Computer Security* abs/1708.06733.
- Guo, J., Fan, Y., Pang, L., Yang, L., Ai, Q., Zamani, H., Wu, C., Croft, W.B., Cheng, X., 2019. A deep look into neural ranking models for information retrieval. *Inf. Process. Manage.* 102067.
- Guri, M., Elovici, Y., 2018. Bridgeware: the air-gap malware. *Commun. ACM* 61 (4), 74–82.
- Haim, N., Vardi, G., Yehudai, G., Shamir, O., Irani, M., 2022. Reconstructing training data from trained neural networks. *arXiv preprint arXiv:2206.07758*.
- Han, D., Wang, Z., Zhong, Y., Chen, W., Yang, J., Lu, S., Shi, X., Yin, X., 2020. Practical traffic-space adversarial attacks on learning-based NIDSs. *arXiv preprint arXiv:2005.07519*.
- Hao, K., 2019. Deepfakes may be a useful tool for spies - mit technology review. <https://www.technologyreview.com/2019/06/14/134940/deepfakes-spies-espionage/>, (Accessed on 06/21/2022).
- Hasegawa, K., Yanagisawa, M., Togawa, N., 2020. Trojan-net classification for gate-level hardware design utilizing boundary net structures. *IEICE Trans. Inf. Syst.* 103 (7), 1618–1622.
- Heuser, A., Picek, S., Guilley, S., Mentens, N., 2016. Side-channel analysis of lightweight ciphers: does lightweight equal easy? In: *International Workshop on Radio Frequency Identification: Security and Privacy Issues*. Springer, pp. 91–104.
- Hidano, S., Murakami, T., Katsumata, S., Kiyomoto, S., Hanaoka, G., 2017. Model inversion attacks for prediction systems: without knowledge of non-sensitive attributes. In: *2017 15th Annual Conference on Privacy, Security and Trust (PST)*, pp. 115–11509. doi:10.1109/PST.2017.00023.
- Hitaj, B., Gasti, P., Ateniese, G., Perez-Cruz, F., 2019. PassGAN: a deep learning approach for password guessing. In: *International Conference on Applied Cryptography and Network Security*. Springer, pp. 217–237.
- Hollnagel, E., Woods, D.D., Leveson, N., 2006. *Resilience Engineering: Concepts and Precepts*. Ashgate Publishing, Ltd.
- Horák, K., Bošanský, B., Tomášek, P., Kiekintveld, C., Kamhoua, C., 2019. Optimizing honeypot strategies against dynamic lateral movement using partially observable stochastic games. *Comput. Secur.* 87, 101579. doi:10.1016/j.cose.2019.101579. <http://www.sciencedirect.com/science/article/pii/S0167404819300665>
- Hussain, M., Al-Haiqi, A., Zaidan, A., Zaidan, B., Kiah, M.M., Anuar, N.B., Abdulaabi, M., 2016. The rise of keyloggers on smartphones: a survey and insight into motion-based tap inference attacks. *Pervasive Mob. Comput.* 25, 1–25.
- Huybrechts, T., Vanommeslaeghe, Y., Blontrock, D., Van Barel, G., Hellinckx, P., 2017. Automatic reverse engineering of can bus data using machine learning techniques. In: *International Conference on P2P, Parallel, Grid, Cloud and Internet Computing*. Springer, pp. 751–761.
- Ilin, I., 2020. Building a news aggregator from scratch: news filtering, classification, grouping in threads and ranking. <https://towardsdatascience.com/building-a-news-aggregator-from-scratch-news-filtering-classification-grouping-in-threads-and-7b0bbf619b68>, (Accessed on 10/14/2020).
- Intelligence, F.E.T., 2015. *HAMMERTOSS: Stealthy Tactics Define a Russian Cyber Threat Group*. FireEye, Inc, Milpitas, CA.
- Ispoglou, K.K., Payer, M., 2016. malWASH: washing malware to evade dynamic analysis. *10th USENIX Workshop on Offensive Technologies (WOOT 16)*. USENIX Association, Austin, TX. <https://www.usenix.org/conference/woot16/workshop-program/presentation/ispoglou>
- Janota, M., 2018. Towards generalization in QBF solving via machine learning. In: *AAAI*, pp. 6607–6614.
- Javed, A.R., Beg, M.O., Asim, M., Baker, T., Al-Bayatti, A.H., 2020. AlphaLogger: detecting motion-based side-channel attack using smartphone keystrokes. *J. Ambient Intell. Humaniz. Comput.* 1–14.
- Jia, H., Choquette-Choo, C.A., Chandrasekaran, V., Papernot, N., 2021. Entangled watermarks as a defense against model extraction. *arXiv preprint arXiv:2002.12200*.
- Jia, Y., Zhang, Y., Weiss, R., Wang, Q., Shen, J., Ren, F., Chen, Z., Nguyen, P., Pang, R., Lopez Moreno, I., Wu, Y., 2018. Transfer learning from speaker verification to multispeaker text-to-speech synthesis. In: Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R. (Eds.), *Advances in Neural Information Processing Systems 31*. Curran Associates, Inc., pp. 4480–4490. <http://papers.nips.cc/paper/7700-transfer-learning-from-speaker-verification-to-multispeaker-text-to-speech-synthesis.pdf>
- Jiang, J., Yu, X., Sun, Y., Zeng, H., 2019. A survey of the software vulnerability discovery using machine learning techniques. In: *International Conference on Artificial Intelligence and Security*. Springer, pp. 308–317.
- Jiang, S., Ye, D., Huang, J., Shang, Y., Zheng, Z., 2020. SmartSteganography: light-weight generative audio steganography model for smart embedding application. *J. Netw. Comput. Appl.* 165, 102689.
- Jiao, L., Zhang, F., Liu, F., Yang, S., Li, L., Feng, Z., Qu, R., 2019. A survey of deep learning-based object detection. *IEEE Access* 7, 128837–128868.
- Juuti, M., Szyller, S., Marchal, S., Asokan, N., 2019. PRADA: protecting against DNN model stealing attacks. In: *2019 IEEE European Symposium on Security and Privacy (EuroS&P)*, pp. 512–527. doi:10.1109/EuroSP.2019.00044.
- Knake, R.K., 2017. *A Cyberattack on the U.S. Power Grid*. Technical Report. Council on Foreign Relations. <http://www.jstor.org/stable/resrep05652>
- Kocher, P., Jaffe, J., Jun, B., 1999. Differential power analysis. In: *Annual International Cryptology Conference*. Springer, pp. 388–397.
- Koh, P.W., Liang, P., 2017. Understanding black-box predictions via influence functions. In: *International Conference on Machine Learning (ICML)*.
- Kolosnjaji, B., Demontis, A., Biggio, B., Maiorca, D., Giacinto, G., Eckert, C., Roli, F., 2018. Adversarial malware binaries: evading deep learning for malware detection in executables. In: *26th European Signal Processing Conf. IEEE, Rome*, pp. 533–537.
- Kong, R., Tong, X., 2020. Dynamic weighted heuristic trust path search algorithm. *IEEE Access* 8, 157382–157390.
- Krebs, B., 2014. Target hackers broke in via HVAC company - krebs on security. <https://krebsonsecurity.com/2014/02/target-hackers-broke-in-via-hvac-company/>, (Accessed on 04/15/2021).
- Kumar, A., Biswas, A., Sanyal, S., 2018. eCommerceGAN: a generative adversarial network for e-commerce. *arXiv preprint arXiv:1801.03244*.
- Kumar, A., Rathore, N., 2016. Improving attribute inference attack using link prediction in online social networks. In: *Recent Advances in Mathematics, Statistics and Computer Science*. World Scientific, pp. 494–503.
- Kuppa, A., Le-Khac, N.-A., 2020. Black box attacks on explainable artificial intelligence(XAI) methods in cyber security. In: *2020 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8. doi:10.1109/IJCNN48605.2020.9206780.
- Kurin, V., Godil, S., Whiteson, S., Catanzaro, B., 2019. Improving sat solver heuristics with graph networks and reinforcement learning. *arXiv preprint arXiv:1909.11830*.
- Lavaud, C., Gerzaguet, R., Gautier, M., Berder, O., Nogues, E., Molton, S., 2021. Whiskering devices: a survey on how side-channels lead to compromised information. *J. Hardware Syst. Secur.* 5 (2), 143–168.
- Leetaru, K., 2019. Deep fakes' greatest threat is surveillance video. <https://www.forbes.com/sites/kalevleetaru/2019/08/26/deep-fakes-greatest-threat-is-surveillance-video/?sh=73c35a6c4550>, (Accessed on 04/15/2021).
- Leong, R., Perez, D., Dean, T., 2019. MESSAGE TAP: who's reading your text messages? <https://www.freeeye.com/blog/threat-research/2019/10/message-tap-who-is-reading-your-text-messages.html>.
- Lerman, L., Bontempi, G., Markowitch, O., 2014. Power analysis attack: an approach based on machine learning. *Int. J. Appl. Cryptogr.* 3 (2), 97–115.
- Lerman, L., Bontempi, G., Taieb, S.B., Markowitch, O., 2013. A time series approach for profiling attack. In: *International Conference on Security, Privacy, and Applied Cryptography Engineering*. Springer, pp. 75–94.
- Leslie, N.O., Harang, R.E., Knachel, L.P., Kott, A., 2019. Statistical models for the number of successful cyber intrusions. *CoRR*. abs/1901.04531.
- Leviathan, Y., Matias, Y., 2018. Google duplex: an AI system for accomplishing real-world tasks over the phone.
- Li, H., Shuai, B., Wang, J., Tang, C., 2015. Protocol reverse engineering using LDA and association analysis. In: *2015 11th International Conference on Computational Intelligence and Security (CIS)*. IEEE, pp. 312–316.
- Li, J., Zhou, L., Li, H., Yan, L., Zhu, H., 2019. Dynamic traffic feature camouflaging via generative adversarial networks. In: *2019 IEEE Conference on Communications and Network Security (CNS)*, pp. 268–276. doi:10.1109/CNS.2019.8802772.
- Li, S., Ma, S., Xue, M., Zhao, B.Z.H., 2022. Deep learning backdoors. In: *Security and Artificial Intelligence*, pp. 313–334. doi:10.1007/978-3-030-98795-4_13.
- Li, Y., Ji, S., Lyu, C., Chen, Y., Chen, J., Gu, Q., Wu, C., Beyah, R., 2020. V-Fuzz: vulnerability prediction-assisted evolutionary fuzzing for binary programs. *IEEE Trans. Cybern.*
- Li, Y., Wang, Y., Wang, Y., Ke, L., Tan, Y.-a., 2020. A feature-vector generative adversarial network for evading PDF malware classifiers. *Inf. Sci.* 523, 38–48.
- Li, Y., Yang, X., Wu, B., Lyu, S., 2019b. Hiding faces in plain sight: disrupting ai face synthesis with adversarial perturbations. *arXiv preprint arXiv:1906.09288*.
- Li, Z., Zou, D., Tang, J., Zhang, Z., Sun, M., Jin, H., 2019. A comparative study of deep learning-based vulnerability detection system. *IEEE Access* 7, 103184–103197.
- Li, Z., Zou, D., Xu, S., Chen, Z., Zhu, Y., Jin, H., 2021. VulDeeLocator: a deep learning-based fine-grained vulnerability detector. *IEEE Trans. Dependable Secure Comput.*
- Li, Z., Zou, D., Xu, S., Ou, X., Jin, H., Wang, S., Deng, Z., Zhong, Y., 2018. VulDeePecker: a deep learning-based system for vulnerability detection. In: *Proceedings 2018 Network and Distributed System Security Symposium* doi:10.14722/ndss.2018.23158.
- Liang, J.H., Oh, C., Mathew, M., Thomas, C., Li, C., Ganesh, V., 2018. Machine learning-based restart policy for CDCL SAT solvers. In: *International Conference on Theory and Applications of Satisfiability Testing*. Springer, pp. 94–110.
- Lim, J., Price, T., Monroe, F., Frahm, J., 2020. Revisiting the threat space for vision-based keystroke inference attacks. In: Bartoli, A., Fusiello, A. (Eds.), *Computer*

- Vision - ECCV 2020 Workshops - Glasgow, UK, August 23–28, 2020, Proceedings, Part V. Springer, pp. 449–461. doi:10.1007/978-3-030-68238-5_33.
- Lin, G., Wen, S., Han, Q.-L., Zhang, J., Xiang, Y., 2020. Software vulnerability detection using deep neural networks: a survey. *Proc. IEEE* 108 (10), 1825–1848.
- Liu, B., Huo, W., Zhang, C., Li, W., Li, F., Piao, A., Zou, W., 2018. α Diff: cross-version binary code similarity detection with DNN. In: Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering, pp. 667–678.
- Liu, H., Lang, B., 2019. Machine learning and deep learning methods for intrusion detection systems: a survey. *Appl. Sci.* 9 (20), 4396.
- Liu, J., Wang, Y., Kar, G., Chen, Y., Yang, J., Gruteser, M., 2015. Snooping keystrokes with mm-level audio ranging on a single phone. In: Proceedings of the 21st Annual International Conference on Mobile Computing and Networking, pp. 142–154.
- Liu, X., Zhou, Z., Diao, W., Li, Z., Zhang, K., 2015. When good becomes evil: keystroke inference with smartwatch. In: Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security, pp. 1273–1285.
- Liu, Y., Ma, S., Aafer, Y., Lee, W.-C., Zhai, J., Wang, W., Zhang, X., 2017. Trojaning attack on neural networks.
- Lowd, D., Meek, C., 2005. Adversarial learning. In: Proc. 11th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD). ACM Press, Chicago, IL, USA, pp. 641–647.
- Lowd, D., Meek, C., 2005. Good word attacks on statistical spam filters. In: Second Conference on Email and Anti-Spam (CEAS). Mountain View, CA, USA.
- Lu, L., Yu, J., Chen, Y., Zhu, Y., Xu, X., Xue, G., Li, M., 2019. KeyLuster: inferring keystrokes on qwerty keyboard of touch screen through acoustic signals. In: IEEE INFOCOM 2019-IEEE Conference on Computer Communications. IEEE, pp. 775–783.
- Lunghi, D., Horejsi, J., Pernet, C., 2017. Untangling the patchwork cyberespionage group. <https://documents.trendmicro.com/assets/tech-brief-untangling-the-patchwork-cyberespionage-group.pdf>.
- Ma, Z., Sheng, O., Pant, G., 2009. Discovering company revenue relations from news: a network approach. *Decis. Support Syst.* 47, 408–414. doi:10.1016/j.dss.2009.04.007.
- Maghrebi, H., Portigliatti, T., Prouff, E., 2016. Breaking cryptographic implementations using deep learning techniques. In: International Conference on Security, Privacy, and Applied Cryptography Engineering. Springer, pp. 3–26.
- Mahadi, N.A., Mohamed, M.A., Mohamad, A.I., Makhtar, M., Kadir, M.F.A., Mamat, M., 2018. A survey of machine learning techniques for behavioral-based biometric user authentication. *Recent Adv. Cryptogr. Netw. Secur.* 43–54.
- Maiorca, D., Demontis, A., Biggio, B., Roli, F., Giacinto, G., 2020. Adversarial detection of flash malware: limitations and open issues. *Comput. Secur.* 96, 101901. doi:10.1016/j.cose.2020.101901.
- Maiti, A., Jadhwal, M., He, J., Bilogrevic, I., 2018. Side-channel inference attacks on mobile keypads using smartwatches. *IEEE Trans. Mob. Comput.* 17 (9), 2180–2194.
- Malhotra, A., Totti, L., Meira Jr, W., Kumaraguru, P., Almeida, V., 2012. Studying user footprints in different online social networks. In: Proceedings of the 2012 International Conference on Advances in Social Networks Analysis and Mining (ASONAM 2012). IEEE Computer Society, USA, pp. 1065–1070. doi:10.1109/ASONAM.2012.184.
- Manning, M., Wong, G.T., Graham, T., Ranbaduge, T., Christen, P., Taylor, K., Wortley, R., Makkai, T., Skorich, P., 2018. Towards a 'smart' cost-benefit tool: using machine learning to predict the costs of criminal justice policy interventions. *Crime Sci.* 7 (1), 12.
- Marquardt, P., Verma, A., Carter, H., Traynor, P., 2011. (sp)iPhone: decoding vibrations from nearby keyboards using mobile phone accelerometers. In: Proceedings of the 18th ACM Conference on Computer and Communications Security, pp. 551–562.
- Martorella, C., 2020. laramies/metagoofil: metadata harvester. <https://github.com/laramies/metagoofil>, (Accessed on 10/20/2020).
- Matta, M., Cardarilli, G.C., Di Nunzio, L., Fazzolari, R., Giardino, D., re, M., Silvestri, F., Spanò, S., 2019. Q-RTS: a real-time swarm intelligence based on multi-agent q-learning. *Electron. Lett.* doi:10.1049/el.2019.0244.
- Mattei, T.A., 2017. Privacy, confidentiality, and security of health care information: lessons from the recent wannacry cyberattack. *World Neurosurg.* 104, 972–974. doi:10.1016/j.wneu.2017.06.104.
- Messaoud, B.I., Guennoun, K., Wahbi, M., Sadik, M., 2016. Advanced persistent threat: new analysis driven by life cycle phases and their challenges. In: 2016 International Conference on Advanced Communication Systems and Information Security (ACOSIS). IEEE, pp. 1–6.
- Mirsky, Y., Doitshman, T., Elovici, Y., Shabtai, A., 2018. Kitsune: an ensemble of autoencoders for online network intrusion detection. 25th Annual Network and Distributed System Security Symposium, NDSS 2018, San Diego, California, USA, February 18–21, 2018. The Internet Society. http://wp.internetsociety.org/ndss/wp-content/uploads/sites/25/2018/02/ndss2018_03A-3_Mirsky_paper.pdf
- Mirsky, Y., Lee, W., 2021. The creation and detection of deepfakes: a survey. *ACM Comput. Surv. (CSUR)* 54 (1), 1–41.
- Mirsky, Y., Mahler, T., Shelef, I., Elovici, Y., 2019. CT-GAN: malicious tampering of 3D medical imagery using deep learning. In: 28th USENIX Security Symposium (USENIX Security 19). USENIX Association, Santa Clara, CA, pp. 461–478. <https://www.usenix.org/conference/usenixsecurity19/presentation/mirsky>
- Mokhov, S.A., Paquet, J., Debbabi, M., 2014. The use of NLP techniques in static code analysis to detect weaknesses and vulnerabilities. In: Sokolova, M., van Beek, P. (Eds.), *Advances in Artificial Intelligence*. Springer International Publishing, Cham, pp. 326–332.
- Monaco, J.V., 2019. What are you searching for? A remote keylogging attack on search engine autocomplete. In: 28th USENIX Security Symposium (USENIX Security 19). USENIX Association, Santa Clara, CA, pp. 959–976. <https://www.usenix.org/conference/usenixsecurity19/presentation/monaco>
- Mozur, P., 2018. Looking through the eyes of China's surveillance state. Accessed: June 2018, <https://www.nytimes.com/2018/07/16/technology/china-surveillance-state.html>.
- Mueller, R., 2018. Indictment - United States of America vs. Viktor Borisovich Netyksho, et al. <https://www.justice.gov/file/1080281/download>.
- Muñoz-González, L., Biggio, B., Demontis, A., Paudice, A., Wonggrassamee, V., Lupo, E.C., Roli, F., 2017. Towards Poisoning of Deep Learning Algorithms with Back-gradient Optimization. In: Thuraisingham, B.M., Biggio, B., Freeman, D.M., Miller, B., Sinha, A. (Eds.), 10th ACM Workshop on Artificial Intelligence and Security. ACM, New York, NY, USA, pp. 27–38.
- Nam, S., Jeon, S., Kim, H., Moon, J., 2020. Recurrent GANs password cracker for IoT password security enhancement. *Sensors* 20 (11), 3106.
- Narayanan, A., Shmatikov, V., 2008. Robust De-anonymization of large sparse datasets. In: 2008 IEEE Symposium on Security and Privacy (sp 2008), pp. 111–125. doi:10.1109/SP.2008.33.
- Nasar, Z., Jaffry, S.W., Malik, M.K., 2019. Textual keyword extraction and summarization: state-of-the-art. *Inform. Process. Manage.* 56 (6), 102088.
- Navalny, A., 2021. European MPs targeted by deepfake video calls imitating Russian opposition - Russia - The Guardian. <https://www.theguardian.com/world/2021/apr/22/european-mps-targeted-by-deepfake-video-calls-imitating-russian-opposition>. (Accessed on 08/17/2022).
- Nicolae, M.-I., Sinn, M., Tran, M.N., Buesser, B., Rawat, A., Wistuba, M., Zantedeschi, V., Baracaldo, N., Chen, B., Ludwig, H., Molloy, I., Edwards, B., 2018. Adversarial robustness toolbox v1.2.0. CoRR. 1807.01069. <https://arxiv.org/pdf/1807.01069>
- Nirkin, Y., Keller, Y., Hassner, T., 2019. FSGAN: subject agnostic face swapping and reenactment. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 7184–7193.
- Novo, C., Morla, R., 2020. Flow-based detection and proxy-based evasion of encrypted malware C2 traffic. In: Proceedings of the 13th ACM Workshop on Artificial Intelligence and Security. Association for Computing Machinery, New York, NY, USA, pp. 83–91. doi:10.1145/3411508.3421379.
- Orekondy, T., Schiele, B., Fritz, M., 2019. Knockoff nets: stealing functionality of black-box models. pp. 4954–4963. https://openaccess.thecvf.com/content_CVPR_2019/html/Orekondy_Knockoff_Nets_Stealing_Functionality_of_Black-Box_Models_CVPR_2019_paper.html.
- Otsuka, H., Watanabe, Y., Matsumoto, Y., 2015. Learning from before and after recovery to detect latent misconfiguration. In: 2015 IEEE 39th Annual Computer Software and Applications Conference, Vol. 3. IEEE, pp. 141–148.
- Ou, X., Govindavajhala, S., Appel, A.W., 2005. MulVAL: a logic-based network security analyzer. USENIX Security Symposium.
- Our work with the DNC: setting the record straight. <https://www.crowdstrike.com/blog/bears-midst-intrusion-democratic-national-committee/>, 2020.
- Oxylabs, 2021. Innovative proxy service to gather data at scale. <https://oxylabs.io/>, (Accessed on 04/14/2021).
- Panagiotou, A., Ghita, B., Shiaeles, S., Bendiab, K., 2019. FaceWallGraph: using machine learning for profiling user behaviour from facebook wall. In: Galinina, O., Andreev, S., Balandin, S., Koucheryavy, Y. (Eds.), *Internet of Things, Smart Spaces, and Next Generation Networks and Systems*. Springer International Publishing, Cham, pp. 125–134.
- Papernot, N., Faghri, F., Carlini, N., Goodfellow, I., Feinman, R., Kurakin, A., Xie, C., Sharma, Y., Brown, T., Roy, A., Matyasko, A., Behzadan, V., Hambardzumyan, K., Zhang, Z., Juang, Y.-L., Li, Z., Sheatsley, R., Garg, A., Uesato, J., Gierke, W., Dong, Y., Berthelot, D., Hendricks, P., Rauber, J., Long, R., 2018. Technical report on the cleverhans v2.1.0 adversarial examples library. arXiv preprint arXiv:1610.00768.
- Pasandi, G., Nazarian, S., Pedram, M., 2019. Approximate logic synthesis: a reinforcement learning-based technology mapping approach. In: 20th International Symposium on Quality Electronic Design (ISQED). IEEE, pp. 26–32.
- Patel, M.I., Suthar, S., Thakar, J., 2019. Survey on image compression using machine learning and deep learning. In: 2019 International Conference on Intelligent Computing and Control Systems (ICCS). IEEE, pp. 1103–1105.
- Peck, J., Nie, C., Sivaguru, R., Grumer, C., Olumofin, F., Yu, B., Nascimento, A., De Cock, M., 2019. CharBot: a simple and effective method for evading DGA classifiers. *IEEE Access* 7, 91759–91771.
- Pellet, H., Shiaeles, S., Stavrou, S., 2019. Localising social network users and profiling their movement. *Comput. Secur.* 81, 49–57.
- Perianin, T., Carré, S., Dyseryn, V., Facon, A., Guillely, S., 2020. End-to-end automated cache-timing attack driven by machine learning. *J. Cryptogr. Eng.* 1–12.
- Perin, G., Chmielewski, Ł., Batina, L., Picek, S., 2021. Keep it unsupervised: horizontal attacks meet deep learning. *IACR Trans. Cryptogr. HardwareEmbedded Syst.* 343–372.
- Picek, S., Heuser, A., Jovic, A., Bhasin, S., Regazzoni, F., 2019. The curse of class imbalance and conflicting metrics with machine learning for side-channel evaluations. *IACR Trans. Cryptogr. HardwareEmbedded Syst.* 2019 (1), 1–29.
- Picek, S., Samiotis, I.P., Kim, J., Heuser, A., Bhasin, S., Legay, A., 2018. On the performance of convolutional neural networks for side-channel analysis. In: International Conference on Security, Privacy, and Applied Cryptography Engineering. Springer, pp. 157–176.
- Pierazzi, F., Pendlebury, F., Cortellazzi, J., Cavallaro, L., 2020. Intriguing properties of adversarial ML attacks in the problem space. In: 2020 IEEE Sympo-

- sium on Security and Privacy (SP), pp. 1332–1349. doi:10.1109/SP40000.2020.00073.
- Pintor, M., Demetrio, L., Sotgiu, A., Melis, M., Demontis, A., Biggio, B., 2022. SecML: secure and explainable machine learning in python. *SoftwareX* 18, 101095. , 2021. Preparing for AI-Enabled Cyberattacks. MIT Technology Review Insights.
- Rahman, T., Rochan, M., Wang, Y., 2019. Video-based person re-identification using refined attention networks. In: 2019 16th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS), pp. 1–8. doi:10.1109/AVSS.2019.8909869.
- Rathi, M., Malik, A., Varshney, D., Sharma, R., Mendiratta, S., 2018. Sentiment analysis of tweets using machine learning approach. In: 2018 Eleventh International Conference on Contemporary Computing (IC3), pp. 1–3. doi:10.1109/IC3.2018.8530517.
- Rebryk, Y., Beliaev, S., 2020. ConVoice: real-time zero-shot voice style transfer with convolutional network. *arXiv preprint arXiv:2005.07815*.
- Ren, Y., Tan, X., Qin, T., Zhao, S., Zhao, Z., Liu, T., 2019. Almost unsupervised text to speech and automatic speech recognition. In: Chaudhuri, K., Salakhutdinov, R. (Eds.), Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9–15 June 2019, Long Beach, California, USA. PMLR, pp. 5410–5419. <http://proceedings.mlr.press/v97/ren19a.html>
- Ribeiro, M.T., Singh, S., Guestrin, C., 2016. “Why should i trust you?”: Explaining the predictions of any classifier. In: 22nd ACM SIGKDD Int'l Conf. Knowl. Disc. Data Mining. ACM, New York, NY, USA, pp. 1135–1144.
- Rigaki, M., Garcia, S., 2018. Bringing a GAN to a knife-fight: adapting malware communication to avoid detection. In: 2018 IEEE Security and Privacy Workshops (SPW), pp. 70–75. doi:10.1109/SPW.2018.00019.
- Roller, S., Dinan, E., Goyal, N., Ju, D., Williamson, M., Liu, Y., Xu, J., Ott, M., Smith, E.M., Boureau, Y., Weston, J., 2021. Recipes for building an open-domain chatbot. In: Merlo, P., Tiedemann, J., Tsarfaty, R. (Eds.), Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, EACL 2021, Online, April 19, - 23, 2021. Association for Computational Linguistics, pp. 300–325. doi:10.18653/v1/2021.eacl-main.24.
- Salminen, J., Jung, S.-G., Jansen, B.J., 2019. The future of data-driven personas: a marriage of online analytics numbers and human attributes. In: ICEIS (1), pp. 608–615.
- Salminen, J., Rao, R.G., Jung, S.-G., Chowdhury, S.A., Jansen, B.J., 2020. Enriching social media personas with personality traits: a deep learning approach using the big five classes. In: International Conference on Human-Computer Interaction. Springer, pp. 101–120.
- Samulowitz, H., Memisevic, R., 2007. Learning to solve QBF. In: AAAI, Vol. 7, pp. 255–260.
- Schreyer, M., Sattarov, T., Reimer, B., Borth, D., 2019. Adversarial learning of deep-fakes in accounting. *1910.03810*.
- Schwartz, J., Kurniawati, H., 2019. Autonomous penetration testing using reinforcement learning. *arXiv preprint arXiv:1905.05965*.
- Seymour, J., Tully, P., 2016. Weaponizing data science for social engineering: automated E2E spear phishing on twitter. *Black Hat USA* 37, 1–39.
- Seymour, J., Tully, P., 2018. Generative models for spear phishing posts on social media. *arXiv preprint arXiv:1802.05196*.
- Shafahi, A., Huang, W.R., Najibi, M., Suciu, O., Studer, C., Dumitras, T., Goldstein, T., 2018. Poison frogs! targeted clean-label poisoning attacks on neural networks. In: Proceedings of the 32nd International Conference on Neural Information Processing Systems. Curran Associates Inc., Red Hook, NY, USA, pp. 6106–6116.
- Shan, S., Wenger, E., Zhang, J., Li, H., Zheng, H., Zhao, B.Y., 2020. Fawkes: protecting privacy against unauthorized deep learning models. In: 29th {USENIX} Security Symposium ({USENIX} Security 20), pp. 1589–1604.
- shaoanlu, 2020. shaoanlu/faceswap-gan: a denoising autoencoder + adversarial losses and attention mechanisms for face swapping. <https://github.com/shaoanlu/faceswap-GAN>, (Accessed on 10/19/2020).
- Sharif, M., Bhagavatula, S., Bauer, L., Reiter, M.K., 2016. Accessorize to a crime: real and stealthy attacks on state-of-the-art face recognition. In: Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security. ACM, pp. 1528–1540.
- Sharon, Y., Berend, D., Liu, Y., Shabtai, A., Elovici, Y., 2021. TANTRA: timing-based adversarial network traffic reshaping attack. *arXiv preprint arXiv:2103.06297*.
- She, D., Krishna, R., Yan, L., Jana, S., Ray, B., 2020. MTFuzz: fuzzing with a multi-task neural network. In: Devanbu, P., Cohen, M.B., Zimmermann, T. (Eds.), ESEC/FSE '20: 28th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, Virtual Event, USA, November 8–13, 2020. ACM, pp. 737–749. doi:10.1145/3368089.3409723.
- She, D., Pei, K., Epstein, D., Yang, J., Ray, B., Jana, S., 2019. NEUZZ: efficient fuzzing with neural program smoothing. In: 2019 IEEE Symposium on Security and Privacy (SP). IEEE, pp. 803–817.
- Shin, E.C.R., Song, D., Moazzezi, R., 2015. Recognizing functions in binaries with neural networks. In: 24th {USENIX} Security Symposium ({USENIX} Security 15), pp. 611–626.
- Shokri, R., Stronati, M., Song, C., Shmatikov, V., 2017. Membership inference attacks against machine learning models. In: 2017 IEEE Symposium on Security and Privacy (SP), pp. 3–18.
- Shumailov, I., Simon, L., Yan, J., Anderson, R., 2019. Hearing your touch: a new acoustic side channel on smartphones. *arXiv preprint arXiv:1903.11137*.
- Shumailov, I., Zhao, Y., Bates, D., Papernot, N., Mullins, R.D., Anderson, R., 2021. Sponge examples: energy-latency attacks on neural networks. In: IEEE European Symposium on Security and Privacy, EuroS&P 2021, Vienna, Austria, September 6–10, 2021. IEEE, pp. 212–231. doi:10.1109/EuroS&P1992.2021.00024.
- Siarohin, A., Lathuilière, S., Tulyakov, S., Ricci, E., Sebe, N., 2019. First order motion model for image animation. In: Conference on Neural Information Processing Systems (NeurIPS).
- Sidi, L., Nadler, A., Shabtai, A., 2020. MaskDGA: an evasion attack against DGA classifiers and adversarial defenses. *IEEE Access* 8, 161580–161592.
- Singh, S., Thakur, H.K., 2020. Survey of various ai chatbots based on technology used. In: 2020 8th International Conference on Reliability, Infocom Technologies and Optimization (Trends and Future Directions) (ICRITO). IEEE, pp. 1074–1079.
- Software, C. P., 2021. 2021 Cyber security report - check point software. <https://www.checkpoint.com/pages/cyber-security-report-2021/>, (Accessed on 06/23/2022).
- Solomon, A., Michaelshvili, M., Bitton, R., Shapira, B., Rokach, L., Puzis, R., Shabtai, A., 2022. Contextual security awareness: a context-based approach for assessing the security awareness of users. *Knowl. Based Syst.* 246, 108709.
- Song, D.X., Wagner, D.A., Tian, X., 2001. Timing analysis of keystrokes and timing attacks on SSH. *USENIX Security Symposium*, Vol. 2001.
- Spiliotopoulos, D., Margaritis, D., Vassilakis, C., 2020. Data-assisted persona construction using social media data. *Big Data Cognit. Comput.* 4 (3), 21.
- Stupp, C., 2020. Fraudsters Used AI to Mimic CEO's Voice in Unusual Cybercrime Case. <https://www.wsj.com/articles/fraudsters-use-ai-to-mimic-ceos-voice-in-unusual-cybercrime-case-11567157402>, (Accessed on 10/14/2020).
- Suciu, O., Coull, S.E., Johns, J., 2019. Exploring adversarial examples in malware detection. In: 2019 IEEE Security and Privacy Workshops (SPW), pp. 8–14. doi:10.1109/SPW.2019.00015.
- Sun, J., Jin, X., Chen, Y., Zhang, J., Zhang, Y., Zhang, R., 2016. Visible: Video-assisted keystroke inference from tablet backside motion. *NDSS*.
- Sun, Q., Tewari, A., Xu, W., Fritz, M., Theobalt, C., Schiele, B., 2018. A hybrid model for identity obfuscation by face replacement. In: Proceedings of the European Conference on Computer Vision (ECCV), pp. 553–569.
- Sutro, A. G., 2020. Machine-learning based evaluation of access control lists to identify anomalies. https://www.tdcommons.org/dpubs_series/2870.
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., Fergus, R., 2014. Intriguing properties of neural networks. In: International Conference on Learning Representations.
- Tariq, N., 2018. Impact of cyberattacks on financial institutions. *J. Internet Bank. Commerce* 23, 1–11.
- Telegram contest, 2020. <https://github.com/IlyaGusev/tgcontest>, (Accessed on 10/14/2020).
- Truong, T.C., Zelinka, I., Senkerik, R., 2019. Neural swarm virus. In: *Swarm, Evolutionary, and Memetic Computing and Fuzzy and Neural Computing*. Springer, pp. 122–134.
- , 2020. The Emergence of Offensive AI. Forrester.
- Ucci, D., Aniello, L., Baldoni, R., 2019. Survey of machine learning techniques for malware analysis. *Comput. Secur.* 81, 123–147.
- Vincent, J., 2022. Binance executive claims scammers made a deepfake of him - the verge. <https://www.theverge.com/2022/8/23/23318053/binance-comms-crypto-chief-deepfake-scam-claim-patrick-hillmann>. (Accessed on 09/07/2022).
- Wang, B., Gong, N.Z., 2018. Stealing hyperparameters in machine learning. In: 2018 IEEE Symposium on Security and Privacy (SP). IEEE, pp. 36–52.
- Wang, D., Neupane, A., Qian, Z., Abu-Ghazaleh, N.B., Krishnamurthy, S.V., Colbert, E.J., Yu, P., 2019. Unveiling your keystrokes: a cache-based side-channel attack on graphics libraries. *NDSS*.
- Wang, Q., Zhao, W., Yang, J., Wu, J., Hu, W., Xing, Q., 2019. DeepTrust: a deep user model of homophily effect for trust prediction. In: 2019 IEEE International Conference on Data Mining (ICDM). IEEE, pp. 618–627.
- Wang, S., Nepal, S., Rudolph, C., Grobler, M., Chen, S., Chen, T., 2020. Backdoor attacks against transfer learning with pre-trained deep learning models. *IEEE Trans. Serv. Comput.* 1. doi:10.1109/TSC.2020.3000900.
- Wang, X., Yamagishi, J., Todisco, M., Delgado, M., Nautsch, A., Evans, N.W.D., Sahidullah, M., Vestman, V., Kinnunen, T., Lee, K.A., Juvola, L., Alku, P., Peng, Y., Hwang, H., Tsao, Y., Wang, H., Maguer, S.L., Becker, M., Ling, Z., 2020. ASvspoof 2019: a large-scale public database of synthesized, converted and replayed speech. *Comput. Speech Lang.* 64, 101114. doi:10.1016/j.csl.2020.101114.
- Wang, Y., Bao, T., Ding, C., Zhu, M., 2017. Face recognition in real-world surveillance videos with deep learning method. In: *Image, Vision and Computing (ICIVC), 2017 2nd International Conference on*. IEEE, pp. 239–243.
- Wang, Y., Cai, W., Gu, T., Shao, W., 2019. Your eyes reveal your secrets: an eye movement based password inference on smartphone. *IEEE Trans. Mob. Comput.*
- Wang, Y., Cai, W., Gu, T., Shao, W., Khalil, I., Xu, X., 2018. GazeRevealer: inferring password using smartphone front camera. In: Proceedings of the 15th EAI International Conference on Mobile and Ubiquitous Systems: Computing, Networking and Services, pp. 254–263.
- Wang, Y., Jia, P., Liu, L., Huang, C., Liu, Z., 2020. A systematic review of fuzzing based on machine learning techniques. *PLoS ONE* 15 (8), e0237749.
- Wang, Y., Yao, Q., Kwok, J.T., Ni, L.M., 2020. Generalizing from a few examples: a survey on few-shot learning. *ACM Comput. Surv.* 53 (3). doi:10.1145/3386252.
- Wang, Y., Zhang, Z., Yao, D.D., Qu, B., Guo, L., 2011. Inferring protocol state machine from network traces: a probabilistic approach. In: *International Conference on Applied Cryptography and Network Security*. Springer, pp. 1–18.
- Weissbart, L., Picek, S., Batina, L., 2019. One trace is all it takes: machine learning-based side-channel attack on EdDSA. In: Bhasin, S., Mendelson, A., Nandi, M. (Eds.), *Security, Privacy, and Applied Cryptography Engineering*. Springer International Publishing, Cham, pp. 86–105.

- White, A., Matthews, A., Snow, K., Monrose, F., 2011. Phonotactic reconstruction of encrypted VoIP conversations: hookt on fon-iks, pp. 3–18. doi:10.1109/SP.2011.34.
- Woh, S., Lee, J., 2018. Game state prediction with ensemble of machine learning techniques. In: 2018 Joint 10th International Conference on Soft Computing and Intelligent Systems (SCIS) and 19th International Symposium on Advanced Intelligent Systems (ISIS), pp. 89–92. doi:10.1109/SCIS-ISIS.2018.00025.
- Woods, D.D., Hollnagel, E., 2006. *Joint Cognitive Systems: Patterns in Cognitive Systems Engineering*. CRC Press.
- Workman, M., 2008. Wisecrackers: a theory-grounded investigation of phishing and pretext social engineering threats to information security. *J. Am. Soc. Inform. Sci. Technol.* 59 (4), 662–674.
- Wu, R., Gong, J., Tong, W., Fan, B., 2021. Network attack path selection and evaluation based on q-learning. *Appl. Sci.* 11 (1). doi:10.3390/app11010285.
- Xu, X., Liu, C., Feng, Q., Yin, H., Song, L., Song, D., 2017. Neural network-based graph embedding for cross-platform binary code similarity detection. In: *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, pp. 363–376.
- Xu, X., Liu, C., Feng, Q., Yin, H., Song, L., Song, D., 2017. Neural network-based graph embedding for cross-platform binary code similarity detection. In: *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security* doi:10.1145/3133956.3134018.
- Yager, R.R., 1984. Approximate reasoning as a basis for rule-based expert systems. *IEEE Trans. Syst. Man Cybern.* SMC-14 (4), 636–643. doi:10.1109/TSMC.1984.6313337.
- Yang, Z., Hu, Z., Dyer, C., Xing, E.P., Berg-Kirkpatrick, T., 2018. Unsupervised text style transfer using language models as discriminators. In: *Advances in Neural Information Processing Systems*, pp. 7287–7298.
- Yao, Y., Li, H., Zheng, H., Zhao, B.Y., 2019. Latent backdoor attacks on deep neural networks. In: *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*. Association for Computing Machinery, New York, NY, USA, pp. 2041–2055. doi:10.1145/3319535.3354209.
- Ye, F., Zhou, S., Venkat, A., Marucs, R., Tatbul, N., Tithi, J. J., Petersen, P., Mattson, T., Kraska, T., Dubey, P., et al., 2020. MISIM: an end-to-end neural code similarity system. *arXiv preprint arXiv:2006.05265*.
- Yousefi, M., Mtetwa, N., Zhang, Y., Tianfield, H., 2018. A reinforcement learning approach for attack graph analysis. In: 2018 17th IEEE International Conference On Trust, Security And Privacy In Computing And Communications/ 12th IEEE International Conference On Big Data Science And Engineering (TrustCom/BigDataSE), pp. 212–217. doi:10.1109/TrustCom/BigDataSE.2018.00041.
- Yu, J., Lu, L., Chen, Y., Zhu, Y., Kong, L., 2019. An indirect eavesdropping attack of keystrokes on touch screen through acoustic sensing. *IEEE Trans. Mob. Comput.*
- Yun, S., Jeong, M., Kim, R., Kang, J., Kim, H.J., 2019. Graph transformer networks. In: Wallach, H.M., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E.B., Garnett, R. (Eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8–14, 2019, Vancouver, BC, Canada*, pp. 11960–11970. <https://proceedings.neurips.cc/paper/2019/hash/9d63484abb477c97640154d40595a3bb-Abstract.html>
- Zelinka, I., Das, S., Sikora, L., Šenkeřík, R., 2018. Swarm virus-next-generation virus and antivirus paradigm? *Swarm Evol. Comput* 43, 207–224.
- Zeng, W., Church, R.L., 2009. Finding shortest paths on real road networks: the case for a*. *Int. J. Geogr. Inf. Sci.* 23 (4), 531–543. doi:10.1080/13658810801949850.
- zerofox, 2020. zerofox-oss/snap_r: a machine learning based social media pentesting tool. https://github.com/zerofox-oss/SNAP_R. (Accessed on 10/21/2020).
- Zhang, H., Chen, H., Song, Z., Boning, D.S., Dhillon, I.S., Hsieh, C., 2019. The limitations of adversarial training and the blind-spot attack. In: 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6–9, 2019. OpenReview.net. <https://openreview.net/forum?id=HyITBhA5tQ>
- Zhang, M., Chen, Y., 2018. Link prediction based on graph neural networks. In: *Advances in Neural Information Processing Systems*, pp. 5165–5175.
- Zhang, W., Lau, R., Liao, S., Kwok, R.-W., 2012. A probabilistic generative model for latent business networks mining. In: *International Conference on Information Systems, ICIS 2012, vol. 2*, pp. 1102–1118.
- Zhang, X., 2018. Analysis of new agent tesla spyware variant. <https://www.fortinet.com/blog/threat-research/analysis-of-new-agent-tesla-spyware-variant.html>.
- Zhang, Y., Meng, J.E., Pratama, M., 2016. Extractive document summarization based on convolutional neural networks. In: *IECON 2016 - 42nd Annual Conference of the IEEE Industrial Electronics Society*, pp. 918–922. doi:10.1109/IECON.2016.7793761.
- Zhiyang, F., Wang, J., Li, B., Wu, S., Zhou, Y., Huang, H., 2019. Evading anti-malware engines with deep reinforcement learning. *IEEE Access* PP, 1–1. doi:10.1109/ACCESS.2019.2908033.
- Zhou, B., Elbadry, M., Gao, R., Ye, F., 2017. BatMapper: acoustic sensing based indoor floor plan construction using smartphones. In: *Proceedings of the 15th Annual International Conference on Mobile Systems, Applications, and Services*, pp. 42–55.
- Zhu, X., Jing, X., You, X., Zhang, X., Zhang, T., 2018. Video-based person re-identification by simultaneously learning intra-video and inter-video distance metrics. *IEEE Trans. Image Process.* 27 (11), 5683–5695. doi:10.1109/TIP.2018.2861366.
- Yisroel Mirsky** (PhD 2018) is a tenure-track lecturer and Zuckerman Faculty Scholar in the Department of Software and Information Systems Engineering at Ben-Gurion University. He is also the head of the Offensive AI Research Lab at BGU. His main research interests include deepfakes, adversarial machine learning, anomaly detection, and intrusion detection. Dr. Mirsky has published his work in some of the best security venues: USENIX, CCS, NDSS, Euro S&P, Black Hat, DEF CON, RSA, CSF, AISEC, etc. His research has also been featured in many well-known media outlets: Popular Science, Scientific American, Wired, The Wall Street Journal, Forbes, and BBC. Some of his works, include the exposure of vulnerabilities in the US 911 emergency services and research into the threat of deepfakes in medical scans, both featured in The Washington Post.
- Ambra Demontis** is an Assistant Professor at the University of Cagliari, Italy. She received her MSc degree (Hons.) in Computer Science and her PhD degree in Electronic Engineering and Computer Science from the University of Cagliari, Italy, respectively, in 2014 and 2018. Her research interests include secure machine learning, kernel methods, and computer security. She co-organizes the AISEC workshop, serves on the program committee of different conferences and workshops, such as IJCAI and DLS, and as a reviewer for several journals, such as TNNLS, TOPS, Machine Learning, and Pattern Recognition. She is a Member of the IEEE and the IAPR.
- Jaidip Kotak** is a Doctorate student in the Information Systems Engineering department at Ben-Gurion University of the Negev (BGU), Israel. His primary areas of interest are computer and network security, general hacking, and application of machine learning in detecting cyber attacks. Jaidip has a MTech with a specialization in Cyber Security & Incident Response and BE in Information Technology from India.
- Ram Shankar** is a data cowboy within the Azure Security Data Science team at Microsoft where he works at the intersection of machine learning and security. He is the founder of the Security Data Science Colloquium – the only avenue where security data scientists from every major cloud provider congregate. Ram is also an affiliate at the Berkman Klein Center at Harvard University.
- Deng Gelei** is a 2nd year PhD student at the School of Computer Science and Engineering, Nanyang Technological University. He received his bachelor's degree from Singapore University of Technology and Design. Before PhD study, he worked as research engineer and penetration tester at Singapore Agency for Science, Technology and Research (A*STAR). His research interests include software security and system security. He has recently been focusing on uncovering vulnerabilities in complex cyber-physical systems through software engineering techniques.
- Liu Yang** graduated in 2005 with a Bachelor of Computing (Honours) in National University of Singapore (NUS). In 2010, he obtained his PhD and started his post-doctoral work in NUS and MIT. In 2012, he joined Nanyang Technological University (NTU), and currently is a full professor and Director of the cybersecurity lab in NTU. Dr. Liu specializes in software engineering, cybersecurity and artificial intelligence. His research has bridged the gap between the theory and practical usage of program analysis, data analysis and AI to evaluate the design and implementation of software for high assurance and security. By now, he has more than 400 publications in top tier conferences and journals. He has received a number of prestigious awards including MSRA Fellowship, TRF Fellowship, Nanyang Assistant Professor, Tan Chin Tuan Fellowship, Nanyang Research Award 2019, ACM Distinguished Speaker, NRF Investigatorship, and 20 best paper awards and one most influential system award in top software engineering conferences like ASE, FSE and ICSE
- Xiangyu Zhang** has substantial research experience in the areas of program analysis, security and AI. His group has produced a list of binary analysis tools that can identify and extract functional components from binaries, reuse binary functional components for rendering forensic evidence in memory, rewrite x86 binaries without relocation information, lift binary execution traces to stand-alone programs that can be compiled and executed on other platforms, vet x86 real world binaries and thousands of iOS apps, and stealthily monitor binary execution. Zhang's group has developed novel attack forensics techniques that instrument applications such that a minimal set of critical application events can be recorded in addition to system level events to facilitate producing comprehensive cyber-attack provenance. Zhang has also substantial experience in AI systems and AI security, especially in AI backdoor attack detection and mitigation.
- Maura Pintor** is a Postdoctoral Researcher at the PRA Lab, in the Department of Electrical and Electronic Engineering of the University of Cagliari, Italy. She received her PhD in Electronic and Computer Engineering from the University of Cagliari in 2022. Her research focuses on machine learning security.
- Wenke Lee** is a Professor and John P. Imlay Jr. Chair in the School of Cybersecurity and Privacy in the College of Computing at The Georgia Institute of Technology. He was also the Director of the Institute for Information Security & Privacy (IISP) at Georgia Tech from 2015 to 2021. He received his PhD in Computer Science from Columbia University in 1999. His research interests include systems and network security, machine learning, and applied cryptography, and he has been the PI of several large-scale projects funded by the government and industry. He is an ACM Fellow and IEEE Fellow.
- Yuval Elovici** is the director of the Telekom Innovation Laboratories at Ben-Gurion University of the Negev (BGU), head of BGU's Cyber Security Research Center, and a professor in the Department of Software and Information Systems Engineering at BGU. He holds BSc and MSc degrees in computer and electrical engineering from

BGU and a PhD in information systems from Tel Aviv University. His primary research interests are computer and network security, cyber security, web intelligence, information warfare, side-channel attacks, AI security, and machine learning. Prof. Elovici also consults professionally in the area of cyber security, sharing his expertise with both startups and large international companies, and is the co-founder of Morphisec, a startup that develops innovative cyber security mechanisms related to moving target defense, and CyberMed which focuses on securing medical imaging devices.

Battista Biggio (MSc 2006, PhD 2010) is Assistant Professor at the University of Cagliari, Italy, and co-founder of the cybersecurity company Pluribus One. He has provided pioneering contributions in machine-learning security, playing a leading role in this field. He has managed six research projects, and regularly serves as a PC member for ICML, NeurIPS, ICLR, IEEE Symp. S&P, and USENIX Security. He chaired IAPR TC1 (2016–2020), co-organized S+SSPR, AISec and DLS, and served as Associate Editor for IEEE TNNLS, IEEE CIM, and Pattern Recognition. He is a senior member of IEEE and ACM, and a member of IAPR and ELLIS.

Israel-US BIRD Foundation