

# Towards Secure MLOps: Surveying Attacks, Mitigation Strategies, and Research Challenges

Raj Patel<sup>\*‡‡</sup>, Himanshu Tripathi<sup>††‡</sup>, Jasper Stone<sup>††‡</sup>, Noorbakhsh Amiri Golilarz<sup>§‡‡</sup>, Sudip Mittal<sup>¶††</sup>, Shahram Rahimi<sup>||‡‡</sup>, and Vini Chaudhary<sup>\*\*††</sup>

<sup>‡‡</sup>Department of Computer Science,

The University of Alabama, Tuscaloosa, AL, USA

<sup>††</sup>Department of Computer Science & Engineering,

Mississippi State University, Mississippi State, MS, USA

raj.patel@ua.edu<sup>\*</sup>, {ht557<sup>†</sup>, jws819<sup>‡</sup>}@msstate.edu, noor.amiri@ua.edu<sup>§</sup> mittal@cse.msstate.edu<sup>¶</sup>,

srahimi1@ua.edu<sup>||</sup>, vchaudhary@cse.msstate.edu<sup>\*\*</sup>

**Abstract**—The rapid adoption of machine learning (ML) technologies has driven organizations across diverse sectors to seek efficient and reliable methods to accelerate model development-to-deployment. Machine Learning Operations (MLOps) has emerged as an integrative approach addressing these requirements by unifying relevant roles and streamlining ML workflows. As the MLOps market continues to grow, securing these pipelines has become increasingly critical. However, the unified nature of MLOps ecosystem introduces vulnerabilities, making them susceptible to adversarial attacks where a single misconfiguration can lead to compromised credentials, severe financial losses, damaged public trust, and the poisoning of training data. Our paper presents a systematic application of the MITRE ATLAS (Adversarial Threat Landscape for Artificial-Intelligence Systems) framework, a comprehensive and continuously updated catalog of AI-focused attacks, to systematically assess attacks across different phases of the MLOps ecosystem. We begin by examining the preparatory phases during which adversaries acquire the essential intelligence required to initiate their attacks. We then present a structured taxonomy of attack techniques explicitly mapped to corresponding phases of the MLOps ecosystem, supported by examples drawn from red-teaming exercises and real-world incidents. This is followed by a taxonomy of mitigation strategies aligned with these attack categories, offering actionable early-stage defenses to strengthen the security of MLOps ecosystem. Given the rapid evolution and adoption of MLOps, we further highlight key research gaps that require immediate attention. Our work emphasizes the importance of implementing robust security protocols from the outset, empowering practitioners to safeguard MLOps ecosystem against evolving cyber attacks.

**Impact Statement**—Machine learning operations (MLOps) streamline the development-to-deployment and maintenance of machine learning (ML) solutions across a variety of industries. However, these systems remain vulnerable to security breaches. Such breaches can disrupt operations, compromise critical services, and erode public trust. Hence, our paper presents a comprehensive strategy that combines established cybersecurity standards with focused attack identification and mitigation strategies. Through systematic analysis and real-world case studies, this roadmap help strengthens ML investments by preserving intellectual property and bolstering user confidence. The insights from this study will help policymakers, business leaders, and the broader workforce to develop better economic strategies, improve regulatory compliance, and promote responsible practices. Overall, our research strengthens the ML ecosystem by offering clear guidelines that enable industries to adopt advanced ML technologies with greater confidence and resilience.

**Index Terms**—Machine Learning Operations (MLOps), Cyber-attack, Security

## I. INTRODUCTION

Machine Learning (ML) technologies have fueled remarkable growth across diverse industries, prompting organizations to swiftly transition models from development into production to secure competitive advantages. To manage this acceleration, companies increasingly adopt the Machine Learning Operations (MLOps) lifecycle framework. First introduced in 2015 [1], MLOps is defined as “a set of practices designed to create an assembly line for building and running machine learning models. It helps companies automate tasks and deploy models quickly, ensuring everyone involved (data scientists, engineers, IT) can cooperate smoothly and monitor and improve models for better accuracy and performance” [2]. MLOps has witnessed increasing adoption as organizations recognize its potential to enhance operational efficiency, support scalable deployments, and ensure reliability within evolving machine learning environments.

The strength of MLOps arises from combining the iterative nature of ML with the Continuous Integration and Continuous Delivery (CI/CD) practices established by DevOps. DevOps refers to “a set of practices for automating the processes between software development and information technology operations teams so that they can build, test, and release software faster and more reliably. The goal is to shorten the systems development life cycle and improve reliability while delivering features, fixes, and updates frequently in close alignment with business objectives” [3]. Integrating DevOps methodologies with ML workflows results in robust and flexible MLOps ecosystem capable of rapidly responding to evolving business and security requirements.

Although MLOps offers substantial benefits, the pressure to accelerate model deployment for competitive gain often results in insufficient attention to security, increasing exposure to vulnerabilities across MLOps ecosystem. With the MLOps market projected to grow by 43% over the next five years [4], protecting these pipelines is increasingly critical. High-profile incidents such as the ShadowRay vulnerability [5],

where insufficient authentication control mechanism allowed attackers to seize compute resources and exfiltrate more than \$1 billion worth of data in September 2023, underline the severity of these risks. Similarly, the Arbitrary Code Execution incident involving Google Colab [6] highlighted the dangers of malicious scripts hidden in shared Jupyter Notebooks, emphasizing the urgent need for rigorous code review processes. Further illustrating these attacks, DeepSeek’s V3 model [7] was allegedly created through repeated unauthorized distillation queries to OpenAI’s models [8]. Additionally, research conducted by SpiderSilk [9] demonstrated how a single misconfiguration compromised critical credentials, leading to unauthorized access to Clearview AI’s sensitive code repositories. Given that Clearview AI’s facial recognition tools is used in matching identities using publicly available images, the implications of such a breach are especially concerning. These examples collectively highlight how security lapses in the MLOps ecosystem can lead to infiltration, data theft, operational disruption, and substantial reputational damage.

Traditional security measures, such as deploying standard malware or virus detection tools within MLOps ecosystem, are insufficient to counter these evolving attacks. Adversaries continuously adapt and refine their tactics, employing defense-aware strategies capable of evading traditional security measures. Skylight’s findings [10] demonstrated how adversarial inputs could evade current malware detection systems. In 2020, Palo Alto Networks’ Security AI research team [11] showcased how generic domain name mutation techniques could defeat Convolutional Neural Network (CNN)-based DGA detection modules. Such mutations evade most ML-based detection systems, underscoring the importance of thorough robustness evaluation before deployment. These attacks often leverage Open Source Intelligence Gathering (OSINT) [12] to maximize infiltration rates while minimizing detection and attribution. Such examples underscore the need to embed security deeply within the MLOps ecosystem, addressing threats at every stage of the ML lifecycle.

To address these challenges, this paper provides the first comprehensive mapping of MLOps-centric attacks using the MITRE ATLAS (Adversarial Threat Landscape for Artificial Intelligence Systems) framework [13], an exhaustive and continuously updated catalog of AI-centric attack lifecycles. Additionally, our analysis incorporates established frameworks such as MITRE ATT&CK [14] and the emerging OWASP (Open Worldwide Application Security Project) Top 10 for Large Language Models [15], identifying key security risks inherent to MLOps ecosystem. By integrating these resources, we enhance our attack assessment and address existing gaps within the MITRE ATLAS framework. Leveraging the ATLAS framework enables organizations to perform structured evaluations of attack vectors, uncover vulnerabilities, and integrate security measures throughout the MLOps ecosystem. Our extensive review of mitigation strategies further assists organizations in protecting mission-critical assets, reinforcing resilience against continuously evolving attacks. The major contributions of this paper are as follows:

- We highlight how adversaries exploit vulnerabilities at each phase of the MLOps ecosystem using the ATLAS

framework to identify potential attack vectors.

- We review a comprehensive taxonomy of MLOps-centric attacks and validate their applicability through examples drawn from real-world incidents and/or red-team exercises documented in the ATLAS framework.
- We present a structured taxonomy of mitigation strategies to enhance the detection and prevention of adversarial attacks on MLOps ecosystem.
- We discuss challenges and provide research recommendations to enhance the security and robustness of the MLOps ecosystem, ensuring resilience against emerging sophisticated cyberattacks.

The remainder of this paper presents a comprehensive analysis of adversarial risks inherent to MLOps ecosystems. Section II introduces the MLOps ecosystem, detailing foundational concepts, operational taxonomy, and real-world implementations. Section III examines pre-exploitation tactics commonly employed by adversaries, emphasizing their role in increasing the likelihood of successful breaches. Section IV leverages the MITRE ATLAS framework to review the evolving attack landscape, illustrating the impact of these threats through empirical evidence from real-world incidents and red-teaming exercises. Section V presents mitigation strategies addressing systemic vulnerabilities within the MLOps ecosystem. Finally, Section VI outlines ongoing research challenges in securing MLOps and identifies opportunities to strengthen defenses through integration of emerging security technologies.

## II. MLOPS OVERVIEW

The concept of Machine Learning Operations (MLOps) emerged in 2015 [1] to bridge the gap between isolated model development to practical deployment. MLOps enhances automation, communication, and monitoring within ML workflow that are inspired by DevOp’s principles, aligning them with organization’s objective and creating tangible value [16].

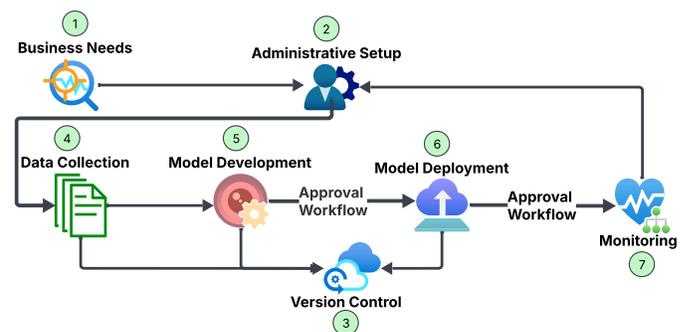


Fig. 1. The MLOps ecosystem begins by identifying business needs, followed by administrative setup, model development, and deployment. Continuous monitoring then provides essential feedback, enabling iterative improvements.

Fig. 1 provides a visual summary of the MLOps ecosystem, which is explored in greater depth in [4]. Here, we outline the key phases most relevant to the security considerations discussed in the following sections. The MLOps process begins with the clear identification of business needs and the establishment of specific Objectives and Key Results (OKRs) that guide decision-making and operational priorities [17].

Effective machine learning models can directly influence commercial outcomes, and a scalable MLOps ecosystem enables the continuous training and deployment of projects as they move from proof-of-concept to production scale [18].

Following business alignment, data collection uses many methods including web scraping, APIs, IoT sensors, and crowdsourcing [16], [18], [19]. Augmentation and transfer learning can enhance the data that is readily available in situations where information is limited [20], [21]. Data preparation typically begins with cleaning procedures that address missing values, errors, and duplicates. This is followed by normalization techniques to ensure data consistency and enhance algorithmic efficiency [19], [22]. Dataset reliability is improved through robust versioning and seamless data integration, which in turn support effective feature engineering and exploratory data analysis (EDA) [18].

A robust administrative framework underpins scalable MLOps infrastructure by supporting key operations such as resource provisioning, user role management, and system monitoring [4]. It ensures access to necessary computational resources, including CPUs, GPUs, and TPUs, for efficient training and deployment. Monitoring tools track system health and anomalies to maintain model consistency and integrity. Additionally, well-defined roles, permissions, and Identity and Access Management (IAM) controls [23] promote secure collaboration and safeguard against unauthorized access.

The model development phase focuses on transforming data into actionable insights through robust experimentation and iterative refinement. Key tasks include data formatting, exploratory data analysis (EDA), feature engineering, and algorithm selection. Using statistical and visual techniques, EDA uncovers patterns, relationships, and anomalies that guide effective feature selection [24], [25]. Feature engineering, through methods such as dimensionality reduction and polynomial expansion, refines raw inputs to improve model accuracy and robustness [18], [20]. Algorithm selection, which includes methods such as decision trees, linear regression, and neural networks, is informed by both data characteristics and organizational needs [17], [20]. Rigorous evaluation metrics, including precision, F1-score, ROC-AUC, and MSE, along with considerations for model interpretability, help ensure alignment with strategic goals and build stakeholder trust [17], [20], [21]. To mitigate overfitting and improve inference, training procedures incorporate techniques such as early stopping, adaptive learning rates, and cross-validation [21], [25].

Version control and experiment tracking allow teams to document and monitor changes to data, features, hyperparameters, and source code, promoting consistency and reproducibility. These systems also provide rollback capabilities in case of performance degradation, helping maintain the stability and reliability of the MLOps workflow [20], [25]. Continuous Integration (CI) automates the testing and validation of new code and model updates, ensuring compatibility and minimizing deployment risks [26], [27].

The deployment phase emphasizes rigorous testing and validation under conditions that closely mirror the production environment, ensuring both reliability and performance. Quality assurance procedures such as batch processing and real-time

inference validations are conducted to verify that models meet predefined accuracy and latency requirements [28]. Responsible AI practices are also incorporated at this stage to mitigate bias and promote fairness. Before release into production, models must pass through a gated approval workflow that includes a combination of automated performance evaluations and manual reviews to ensure compliance with robustness and security standards [4]. Once approved, models are deployed for end-user access and configured to support diverse inference scenarios. Benchmarking tools, including MLPerf [29], are used to assess production readiness and verify that inference quality is maintained in the target environment [30].

Continuous monitoring is essential in complex multi-cloud environments, enabling performance tracking across testing, staging, and production stages to ensure timely detection and resolution of deviations [31]. This comprehensive approach facilitates sustained alignment with evolving business needs.

In summary, the comprehensive MLOps ecosystem effectively manages the machine learning lifecycle, integrating continuous development, deployment, and monitoring. By prioritizing alignment with business objectives and responsible AI practices, organizations can rigorously evaluate and continuously monitor the performance of their machine learning systems. These efforts contribute to operational efficiency, ethical model behavior, sustained reliability, and overall trustworthiness. However, despite these advancements, MLOps ecosystem remain vulnerable to security threats. Ensuring long-term dependability and stakeholder trust requires ongoing research and integration of robust, context-aware security.

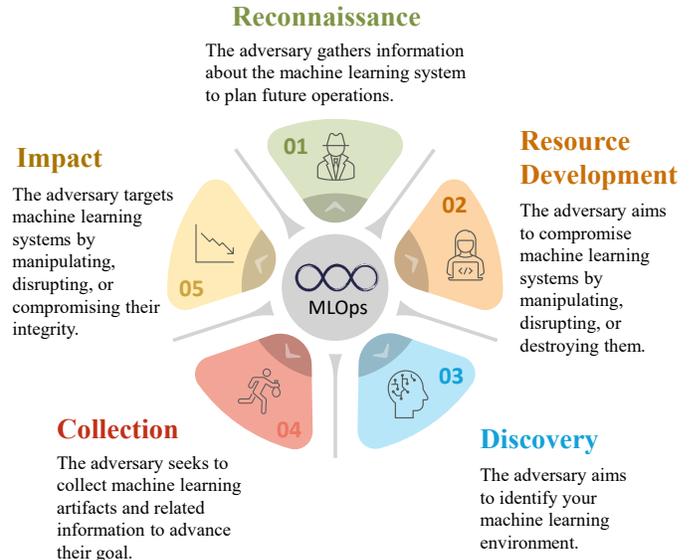


Fig. 2. Five key MITRE ATLAS tactics used throughout an attack lifecycle, beginning chronologically with Reconnaissance (1) and followed by Resource Development (2), Discovery (3), Collection (4), and Impact (5).

### III. SURVEY APPROACH AND TAXONOMY

In the previous section, we examined the MLOps ecosystem, outlining the integration of tools, roles, and software components necessary to support a high degree of automation. While this integration streamlines the ML workflow, it also introduces vulnerabilities at various stages, potentially exposing

the MLOps ecosystem to exploitation and compromising its integrity. In this paper, we present a structured mapping of adversarial tactics and techniques as defined by the MITRE ATLAS framework to stages of the MLOps lifecycle [13]. In this context, a tactic refers to the strategic objective an adversary aims to achieve, while a technique describes the specific method used to accomplish that goal [32], [33]. The MITRE ATLAS framework was selected for its specialized focus on adversarial threats targeting AI systems. It functions as an extension of the broader MITRE ATT&CK framework, which provides comprehensive coverage of cybersecurity threats across domains [14]. Drawing on ATLAS, we conduct a systematic review of how adversarial techniques correspond to vulnerabilities within the MLOps ecosystem. To the best of our knowledge, this is the first study to present a detailed alignment between the MITRE ATLAS framework and stages of the MLOps lifecycle.

Before examining MLOps-specific attacks and techniques, it is important to first consider the foundational phases that enable their success. Prior to exploiting vulnerabilities within an organization, adversaries typically begin with open-source intelligence gathering (OSINT) [12], along with several other preparatory phases that set the stage for later compromise. Fig. 2 highlight key adversarial tactics commonly employed during the attack lifecycle, presented in a typical chronological sequence. Fig. 3 provides a more detailed breakdown of specific techniques organized under each of these tactics. In the following subsections, we will explore the tactics outlined in the MITRE ATLAS framework, focusing on those commonly used in the early stages of an attack, which we refer to as the preparatory or common phases. We will also examine how these tactics enable lateral movement and escalate into more destructive forms of attack.

1) *Reconnaissance (ID: AML.TA0002)*: Reconnaissance serves as the initial phase of an attack, during which adversaries attempt to gather information about the victim organization’s ML capabilities and ongoing research efforts [34]. It begins by scouring publicly available information on the victim organizations. To showcase their ML capabilities, these organizations highlight their ML research efforts by publishing their work in conferences or journal proceedings [35]. These publications highlight their efforts on how they build and deploy these ML models [35]. Similarly, some organizations would publish their un-reviewed research work into platforms like arXiv to get ahead of the competition before the final literature materializes for publication [36]. The Pre-Print repositories are completely open to public at large. In some cases, conference and journal proceedings are hidden behind membership dues but the adversary do not have to worry about that with pre-print articles.

Another valuable source of information are technical blogs where researchers would like to boast about their R&D efforts in the space of machine learning as they typically outline more specific information about practical implementation of their efforts [37]. An adversary also gathers information on common adversarial vulnerabilities in popular ML models. This helps them adapt existing attack methods or develop new ones [38]. Victim’s websites is another valuable resource as it often

details about ML products, services, departments, divisions as well as employee’s involvement [39]. Additionally, attackers can craft search queries to identify ML-capable applications in stores such as Google Play, Apple Store and Microsoft Store [40]. In addition, active system scanning (e.g., port and vulnerability scans) helps identify weaknesses that passive research might miss [41].

Hence, reconnaissance is recognized as a critical first step in the attack phase. Several of these techniques have been observed in scenarios where attackers reverse-engineer machine learning models to conduct model-stealing attacks [42]. For instance, notable cases include the 2020 bypass of Cylance’s AI malware detection system [10] and the evasion of CNN-based DGA detection systems [11]. These incidents demonstrate how reconnaissance equips attackers with the knowledge necessary to launch highly targeted attacks, underscoring the need for robust security throughout the MLOps ecosystem. Following reconnaissance, adversaries typically transition to the resource development phase, where they acquire the tools and capabilities required to exploit identified vulnerabilities.

2) *Resource Development (ID: AML.TA0003)*: In the resource development phase, adversaries acquire or create assets to support subsequent attack operations. These assets may include ML artifacts, rented or purchased infrastructure, established accounts, and specialized tools that enable various stages of the attack lifecycle, such as ML attack staging [43]. After the reconnaissance phase, adversaries will search for public ML artifacts (e.g., software stacks, training data, or model parameters) by scouring cloud storage, public services, and data repositories. They particularly seek ML artifacts that are linked to the victim organization, as these can reveal valuable insights about the victim’s ML setups. In some cases, adversaries register for accounts to gain access to these artifacts, enabling them to craft and inject adversarial inputs to help prepare a proxy or partial production models [44]. Public datasets similar to those used by the target may also be collected to tailor attacks more effectively [45], along with models in formats such as ONNX, HDF5, Pickle, PyTorch, or TensorFlow that align with the victim’s infrastructure [46].

Another tactic involves obtaining specialized or general software capabilities, particularly tools tailored for ML-based attacks. Adversaries may leverage open-source adversarial ML frameworks such as CleverHans [47], the Adversarial Robustness Toolbox (ART) [48], or FoolBox [49], as well as repurpose general-purpose tools for malicious use [50]. In addition, legitimate tools may be customized to compromise ML systems, even when they are not inherently ML-focused. In more advanced cases, adversaries develop their own capabilities, such as crafting adversarial websites or embedding data-exfiltration code in Jupyter notebooks [51], enabling covert execution of ML-relevant attacks [52].

To bolster these operations, adversaries often acquire or rent infrastructure, including servers, domains, mobile devices, and cloud computing resources [53]. They may also develop physical countermeasures, such as adversarial prints or wearable devices, to interfere with sensor input and degrade model performance [54]. For attack staging and experimentation, powerful GPUs and development workspaces are frequently

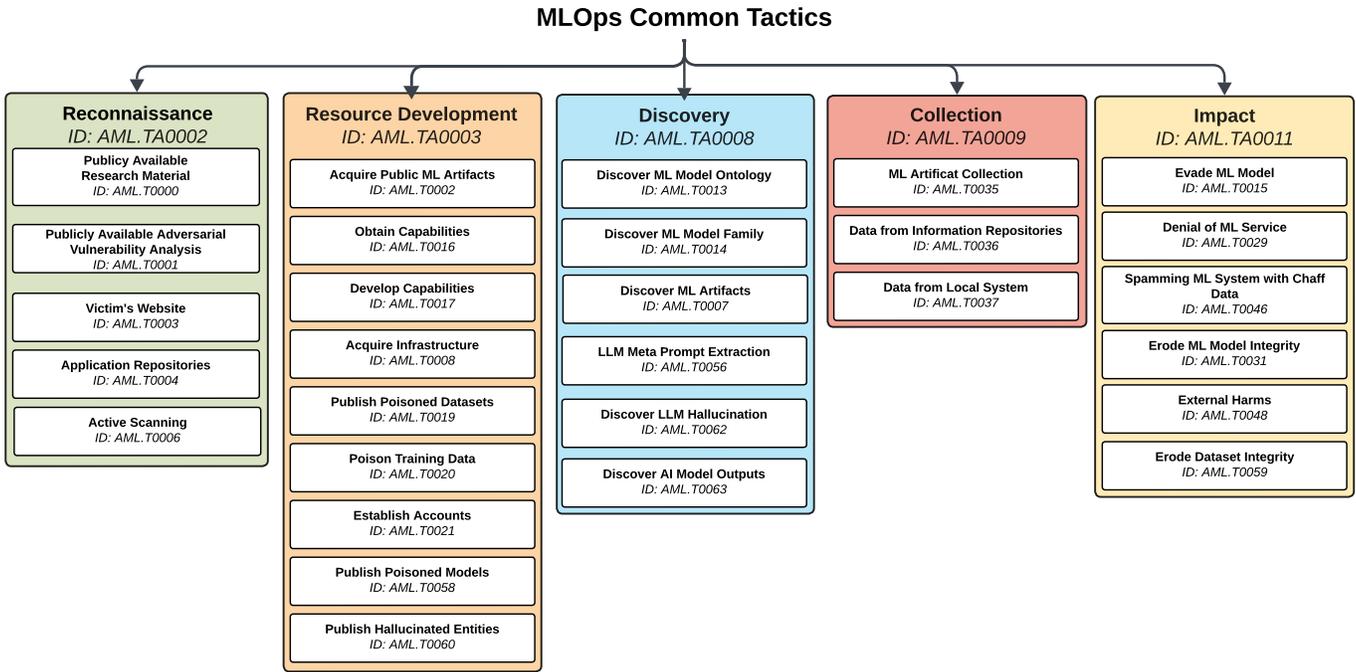


Fig. 3. Expanded overview of the MITRE ATLAS tactics: Reconnaissance, Resource Development, Discovery, Collection, and Impact, mapped to their corresponding techniques. Each tactic is shown as a column anchor, illustrating methods adversaries may use across the MLOps attack lifecycle. Clustering techniques under respective tactics highlights interdependencies and offers a structured view for identifying and countering adversarial behaviors.

required. These can be accessed through free platforms like Google Colab [55] or commercial cloud providers such as AWS [56], Azure [57], or Google Cloud. Distributing operations across multiple environments helps obscure malicious activity by blending it with normal network traffic [58]. Adversaries may register for new domains to host poisoned datasets or exploit expired domains in public repositories by injecting malicious content into previously trusted URLs [59].

In addition to infrastructure, adversaries weaponize data sources by publishing poisoned datasets in public repositories, compromising any models trained with them through ML Supply Chain Compromise [60]. They can also poison training data by modifying inputs or labels, embedding hidden vulnerabilities activated when specific triggers appear [61]. Similarly, they may publish poisoned models in public registries. These are either newly developed or modified to be introduced into victim environments via the supply chain compromise [62]. Furthermore, adversaries create or publish malicious entities (e.g., software packages, websites, email addresses) that might appear in large language model (LLM) generated text, undermining trust in LLM outputs [63]. Lastly, they can create accounts to gain access or impersonate victims, as adversaries establish accounts with various services to further ML attack staging [64].

In the broader context of ML attacks, resource development involves proactively gathering, creating, or altering ML artifacts, infrastructure, and accounts to form a solid base for targeted attacks. Unlike passive reconnaissance, resource development actively prepares the tools and capabilities that power subsequent phases, such as discovery and exploitation. Section IV provides a more detailed explanation of how these resources are ultimately leveraged in the attack phase.

3) *Discovery (ID: AML.TA0008)*: After an adversary has created a foothold in the victim's organization. They look to employ discovery phase to further collect information gathering post-compromise. In the discovery phase, adversaries aim to learn more about victim's ML environment, including its system, environment, and internal network. By gathering this information using native operating system tools, adversary can determine which elements of the environment they can manipulate to advance their objectives [65].

An adversaries may begin by searching for ML artifacts such as software stacks, model zoos, and data management systems to map how models are trained, deployed, and stored [66]. Once these are identified, they move on to discover the model's ontology. In order to discover model's ontology, they may repeatedly query the model or review its documentation to uncover the model's output space. This knowledge reveals how the model is used and enables adversaries to design attacks that exploit the model's specific capabilities [67]. They further try to discover the general family of model it belongs to by analyzing responses or documentation. This information can uncover specific (or common) exploitation methods associated with that model family, which attackers may leverage in crafting their own exploits [68].

In the case of LLMs, adversaries may attempt to extract meta prompts to discover initial configuration instructions that govern the internal workings of the model. They try to extract this information to steal intellectual property of the model [69]. They may further try to take advantage of the inherent trust of the users in hallucinated outputs such as fictitious commands, packages, or URLs to use to exploit and compromise end-users [70]. Some AI model output may contain additional information such as class scores, logs, API responses, etc.

that can allow adversaries to identify weakness and craft attacks based on that information [71]. The above claim was supported by the MITRE AI Red Team exercise in 2020 [72] where they discovered model’s ontology of a commercial facial identification service through its inference API. This allowed them to developed adversarial inputs that led to misclassification. Discovery helps identify internal working of the model to develop a more tailored attack. After fully exploring the environment, adversaries often move to the collection phase, where they extract data and artifacts, leveraging insights gained in discovery phase.

4) *Collection (ID: AML.TA0009)*: In collection phase, adversaries gather machine learning artifacts and other relevant information that can help achieve their objective from sources like software repositories, container registries, model repositories, and object stores [73]. By collecting information on ML artifact, they uncover information about models, datasets, and telemetry outputs that can be leveraged for theft or future operations [73], [74]. Beyond ML artifacts, attackers often focus on collaborative information repositories such as SharePoint, Confluence, SQL Server as they store a wide range of useful data for advancing their objectives [75]. Additionally, adversaries may search local systems for items such as configuration files, SSH keys, or other sensitive assets to facilitate exfiltration or to further their goal [76].

In 2020, researchers from SpiderSilk [9] discovered a misconfiguration that gave them unauthorized access to Clearview AI’s private repositories. This allowed the researcher to gain access to stored credentials, secret keys, and copies of production-ready app. Similarly, a Microsoft Azure Service Disruption in 2020 revealed how collection-related vulnerabilities can be exploited; a red team exercise accessed internal Azure model files and training data, then used them to craft adversarial examples for online attacks [77]. These incidents highlight how successful collection phase paves the way for significant disruptions to an ML system’s integrity and confidentiality. Once attackers have gathered the necessary artifacts, they often transition to the impact phase, using stolen data and insights to compromise availability, damage operations, or steal intellectual property.

5) *Impact (ID: AML.TA0011)*: The Impact tactic represents the potential consequences an adversarial attack can have on the victim organization. In this phase, an adversary tries to manipulate ML systems to degrade performance, undermine trust, or mask malicious actions, often affecting both business and operational processes. By altering or destroying data, attackers may create a false sense of normalcy while advancing their objectives [78]. One of the technique involves evading ML models through adversarial inputs [79]. This prevents correct attack detection by virus scanners or intrusion detection systems, allowing traditional cyber attacks to proceed undetected [79]. This was denoted in the real-world incident in 2020 where attack on camera hijack was carried out on facial recognition system in China to bypass verification systems to allow fraudulent invoices to be issued [80]. Adversaries may also initiate denial of ML Service by overwhelming systems with resource-intensive requests [81], or deploy chaff data to burden analysts with false positives output [82]. Over time,

these strategies erode confidence in the model’s reliability and force costly manual reviews [83].

Financial costs are another target, as adversaries exploit the high computational demands of ML systems for cost-harvesting attacks [84]. These attacks causes victim organization to suffer from various types of external harms such as financial, reputational, societal and user-level damage [85]. Intellectual property theft remains a serious attack as well. By exfiltrating proprietary models or datasets, attackers undermine an organization’s competitive advantage and Machine-Learning-as-a-Service (MLaaS) revenues [86]. Dataset integrity can also suffer when adversaries poison or corrupt data, reducing reliability and forcing organizations to invest resources in remediation [87].

A case study from 2021 highlighted a backdoor attack on deep learning models using “neural payload injection”, triggered by a visual cue [88]. Additionally, a 2020 incident involving ID.me demonstrated how stolen identities were exploited to file fraudulent claims [89]. These incidents showcase how adversaries progressively advance from initial reconnaissance to significantly impacting ML environments. The following subsection now shift focus to elaborate specifically on how adversaries exploit user-level and system-level vulnerabilities within the MLOps ecosystem to undermine model security and disrupt operational integrity.

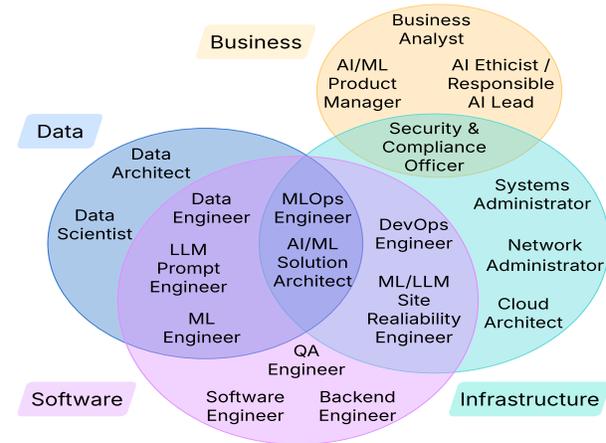


Fig. 4. Venn diagram of potential stakeholders interested in MLOps security.

#### A. User and System-Level

In the context of MLOps ecosystem security, *User-Level* and *System-Level* attacks represent two distinct attack vectors that adversaries exploit across the MLOps lifecycle. As shown in Fig. 4, the stakeholders align with user level roles and responsibilities that are central to the discussion in this subsection. **User-level** attacks in MLOps take aim at the human element, relying on social engineering, credential theft, and deceptive tactics to breach defenses. Adversaries often use phishing and spearphishing [90]–[93] to gather personal or sensitive data from unsuspecting users, sometimes by distributing fraudulent emails or impersonating trusted colleagues. Another method involves inserting “hallucinated entities” [63] into model outputs, prompting individuals to click malicious links or download compromised files. Attackers also gain unauthorized

access by stealing or abusing user credentials [94], which can grant them control over inference APIs or let them modify model performance. These intrusions are primarily user-centric because they exploit personal trust, habits, and human error rather than purely technical flaws.

By contrast, **System-Level** attacks focus on critical infrastructure integral to MLOps, such as computational resources, pipelines, and version control systems. Attackers may exploit vulnerabilities in GPU or TPU drivers [95], misuse compute power, or alter code repositories that store ML artifacts. Data poisoning is also common at this level [96], where adversaries inject deceptive samples or manipulate labels [22], [97] to undermine model integrity. Some attackers modify the models themselves by embedding backdoors [98] or developing proxy models offline [99]. During runtime, they may exfiltrate private data through inference APIs [79] or craft adversarial inputs that bypass traditional defenses. System-level attacks can arise anywhere in the MLOps ecosystem, including third-party libraries, pre-trained models, or unsafe artifacts. Attackers can also disrupt availability and reliability by launching denial-of-service [81], spoofing [100], or session hijacking [101] attacks.

Overall, user-level exploits tend to be more subtle and hinge on manipulating individual actions and trust, while system-level attacks target the foundational components of MLOps infrastructure. Countering user-level risks involves stronger security policies, user education, and careful management of credentials. On the other hand, system-level defenses require infrastructure hardening, supply chain security, continuous data validation, and network safeguards. Each category poses unique challenges and demands a targeted approach. The next section reviews the attack landscape in greater detail, supported by real-world incidents and red-teaming exercises that demonstrate how adversaries successfully exploit vulnerabilities across the MLOps ecosystem.

#### IV. REVIEW ATTACK AND CHALLENGES

In the preceding section, we outlined the foundational prerequisites or preparatory phases, along with the potential impact an adversarial attack can have on the victim organization. These steps often serve as precursors to deploying exploits and establishing a persistent foothold within the victim organization. In this section, we examine a range of attacks relevant to the MLOps ecosystem, categorizing them into two primary types: (1) those targeting the users involved in the creation and deployment of MLOps systems, and (2) those targeting the systems themselves. As shown in Fig. 5, each category is further subdivided into several subcategories, enabling a structured review of the attack landscape. These subcategories can all be interpreted as variations of supply chain attacks, whether directed at users or systems. Each subsection presents a focused analysis of attacks within its scope. This includes attacks documented in the MITRE ATLAS framework, as well as those not currently included but deemed critical for inclusion in future iterations. To substantiate our review, we support each attack with evidence from real-world incidents or red-teaming exercises, highlighting their practical effectiveness and relevance. Together, this comprehensive review highlights

the complexity of the MLOps threat landscape and reinforces the need for robust, multi-layered security strategies.

##### A. Exploitation Attacks

Humans are often considered the weakest link in cybersecurity [102]. Exploitation attacks target vulnerabilities among users within the MLOps ecosystem. These attacks often prompt unintentional user actions that trigger malicious code or packages, ultimately disrupting workflows and compromising operations. They undermine operations through phishing [90] [91], social engineering [93], and privilege escalation [103]. Phishing [90], including spearphishing [91], uses fraudulent messages to gain unauthorized access. Advances in generative AI and deepfake technologies [93] [92] have made phishing more convincing, as attackers craft realistic messages and impersonate trusted individuals to steal credentials.

User execution [104] relies on user actions to activate malicious code introduced through ML Supply Chain Compromise [105] or social engineering [93]. Adversaries trick users into executing harmful files or links, inadvertently granting access. Unsafe ML artifacts [106] embed harmful payloads and may establish persistent access when models are stored, transferred, or loaded. Serialization is common for model, but insufficient checks present an opening for code execution [15]. Malicious packages [107] can masquerade as legitimate AI tools yet deliver harmful payloads upon import. In "ChatGPT Package Hallucination" incident, attackers registered fabricated packages, causing unsuspecting users to install them [108].

Privilege escalation [103] grants attackers elevated privileges. Elevated privileges are often necessary to achieve an adversary's objective. Privilege escalation can be achieved through techniques such as LLM prompt injection (direct [109] and indirect [110]), LLM plugin compromise [111], LLM jailbreak [112], and excessive agency [15]. LLM prompt injection [113] manipulates a model's behavior, while plugin compromise [111] exploits integrated plugins to retrieve sensitive data. Jailbreaking [112] bypasses model restrictions, and excessive agency [15] arises when systems have permissions beyond their scope. The Morris II worm [114] used indirect prompt injection in a retrieval-augmented generation (RAG) email assistant, exfiltrating data and replicating malicious behavior. The ChatGPT plugin privacy leak [115] illustrated how excessive permissions lead to data breaches.

These exploitation attacks underscore the need for robust defenses within the MLOps ecosystem. By leveraging phishing [90] [91] [92], spearphishing [93], privilege escalation [103], and user execution [106] [107] [15], adversaries can disrupt workflows and undermine trust. The next section explores access abuse attacks, which focus on exploiting system accounts and permissions to infiltrate MLOps infrastructure.

##### B. Access Abuse Attacks

Access abuse attacks in the MLOps ecosystem occur when adversaries leverage legitimate credentials or authentication tokens to gain unauthorized access to ML systems. These credentials, such as inference API keys or service tokens, allow attackers to bypass traditional security controls and manipulate

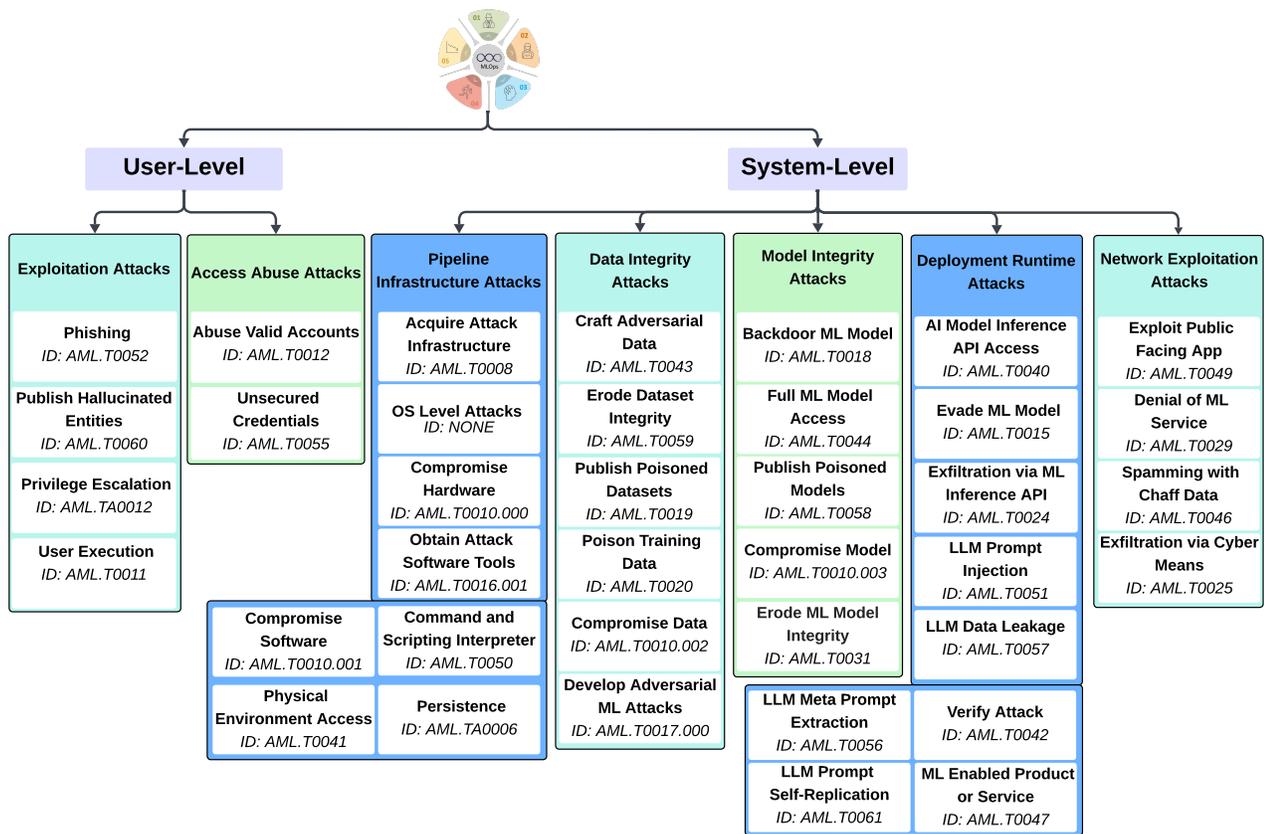


Fig. 5. Attacks are categorized into two main groups: user-level and system-level. The clusters named underneath them to show conceptual alignment with MLOps phases rather than a strict correlation. Some techniques do not have a MITRE ATLAS ID, indicating they originate from other research sources. This structured layout shows how adversaries can employ multiple tactics across the MLOps supply chain.

or disrupt the broader operational environment [116]. Such intrusions can compromise model integrity, expose sensitive data, and interfere with operational workflows, underscoring the importance of robust access control and monitoring.

A common technique involves exploiting valid accounts [94] [117], including default accounts [118], domain accounts [119], local accounts [120], and cloud-based accounts [119]. Default accounts, often overlooked after deployment, and inactive accounts, such as those belonging to former employees, are particularly vulnerable due to inadequate oversight [117]. Cloud-based accounts, including API credentials for platforms such as AWS [56], Azure [57], or Google Colab [55], present an expanded attack surface. Real-world incidents illustrate this risk. In July 2022, a Google Colab vulnerability [6] allowed arbitrary code execution, and a 2020 Azure red team exercise [77] showed how valid accounts were used to infiltrate internal services, exfiltrate data, and escalate privileges.

Another common method involves the use of unsecured credentials [121] [122], which may be stored in plaintext or exposed through environment variables, configuration repositories, and container platforms. Attackers seek sensitive assets such as API keys, tokens, and passwords by exploiting access to the Cloud Instance Metadata API [123], Kubernetes [124], Docker [125], and container logs [126]. Unsecured credentials are also frequently extracted from the Windows Registry [127], bash history files [128], and private keys stored in predictable locations [129]. Automated tools and scripts further streamline this process, broadening the scope of attacks. In addition,

collaboration platforms such as Slack [130], Trello [131], and Microsoft Teams [132] may unintentionally store credentials that attackers can exploit. These factors highlight how unmanaged credentials jeopardize ML operations. High-profile incidents include the MathGPT breach in January 2023 [133], which exposed unsecured GPT-3 API keys, and a Ray framework vulnerability in September 2023 [5], where attackers accessed SSH keys and cloud credentials for data exfiltration and resource abuse.

In conclusion, access abuse-level attacks highlight how weak credential management and authentication can compromise security across the MLOps ecosystem. Addressing these vulnerabilities is essential for protecting the integrity of machine learning operations. The following section explores system-level attacks targeting the infrastructure that supports ML workflows, beginning with pipeline infrastructure attacks.

### C. Pipeline Infrastructure Attacks

Pipeline infrastructure attacks focus on the core systems that support machine learning workflows, including servers, CI/CD pipelines, GPUs, and version control systems as shown in Fig. 1. By compromising these components, adversaries can alter or delay model training, testing, and deployment. Many MLOps ecosystem stages depend on interlinked resources, so disruptions often spread from one area to the next. Malicious actions are frequently disguised within routine operations [134], complicating detection in environments that

rely on continuous delivery. One frequent method is acquiring infrastructure [53] in ways that appear legitimate. Attackers may rent servers [53], register for expired domains [59], or rely on consumer hardware [135], then replicate real workflows to develop targeted exploits. They can also set up cloud-based workspaces [58], such as Google Colab [55], to refine malicious scripts or poison datasets. In some cases, CloudBorne [15] and CloudJacking [15] attacks exploit virtualization layers to gain elevated access and control shared resources. The “split-view” poisoning attack [96] similarly modifies domains embedded in AI datasets, degrading performance or skewing predictions in dependent models.

Hardware compromises [136] also present an additional attack vector within the MLOps ecosystem. GPUs, firmware, and other specialized hardware components can be tampered with at various stages of the supply chain, enabling adversaries to inject code that alters outputs, degrades performance, or exposes model states. Even minor modifications may remain undetected without regular integrity checks. The “LeftOvers” attack (CVE-2023-4969) [15] demonstrated how GPU memory leaks could be exploited to recover sensitive data. Similarly, side-channel attacks [15] allow adversaries to infer proprietary model parameters by observing execution timing or power consumption patterns. When such hardware-level threats go unnoticed, compromised components may persist across multiple training and deployment cycles. Physical access attacks [137] can also introduce vulnerabilities by tampering with sensors or embedding adversarial objects [98] at the data source, potentially leading to systemic model failures.

Another avenue of compromise involves command and scripting interpreters [138]. Python [139], PowerShell [140], and Unix shells are widely used to automate model creation, testing, and deployment. Attackers can insert malicious code appearing as normal updates, leading to adversarial examples [15], corrupted libraries, or altered hyperparameters. As these scripts support CI workflows, such changes may accumulate over time [15]. Memory-related exploits, including buffer overflows [141] or speculative execution attacks [142], can expose or manipulate data across pipeline. Software supply chain attacks [143] offer another pathway, as adversaries may misuse tools [144] or dependencies to introduce malicious packages or backdoored libraries. For example, a dependency confusion attack on PyTorch [145] exfiltrated user data via a package crafted to resemble a trusted component.

Persistence [146] allows attackers to remain in systems for extended periods. Techniques include injecting poisoned training data [61] or distributing backdoored models [98], which appear legitimate but contain triggers or hidden behaviors that activate under specific conditions. Because the MLOps ecosystem often reuses existing assets, these backdoors can remain effective through updates and maintenance. If undetected, they can continue to influence results, hinder collaboration, or exfiltrate data. As these effects spread, they become harder to trace, especially if a compromised GPU enables further tampering in CI scripts. Attackers exploit this complexity [53], bypassing traditional security controls. Comprehensive monitoring, code reviews, and layered defenses help reduce exposure. However, no single strategy mitigates all vulnerabilities across MLOps.

Measures such as infrastructure scanning, code signing, and team training on best practices all play a critical role.

In summary, pipeline infrastructure-level attacks compromise the foundational components of machine learning systems. Hardware tampering [136], script abuse [138], memory exploits [141] [142], and malicious software dependencies [143] each introduce distinct points of entry. Physical access [137] and persistence tactics [146] further amplify these risks. Because such attacks span both technology and operations, a coordinated defense strategy is essential. The discussion now turns to data integrity-level attacks, which focus on manipulating the data that powers machine learning systems.

#### D. Data Integrity Attacks

Data integrity is essential for security and reliability in ML systems. It ensures that training and operational data remain accurate and trustworthy. However, it also creates opportunities for adversaries, who can exploit vulnerabilities in dataset creation, labeling, and deployment processes to corrupt model behavior or mislead outcomes. Poisoning attacks [61] embed harmful data into training sets, thereby compromising model performance. Evasion attacks [97] target models with adversarial inputs during inference, prompting incorrect predictions. Adversaries may also insert false labels or malware into open-source datasets and manipulate non-control data [147], leading to security breaches without altering core system logic.

A central threat is data poisoning, which alters pre-training or fine-tuning data to degrade performance or embed hidden backdoors [15]. Techniques like split-view poisoning [96] exploit gaps in dataset creation by injecting malicious content at specific times. Reliance on external data sources magnifies these risks, as unsafe inputs may enter models lacking robust validation or access controls. In large-scale systems built from dynamic web content [96], adversaries exploit predictable snapshot schedules, embedding attacks shortly before data collection. These threats underscore the need for rigorous data checks and continuous monitoring. Beyond poisoning, adversaries also employ training-only attacks [148], which manipulate training data to degrade model performance or introduce targeted vulnerabilities. Techniques such as feature collision, label flipping, and bilevel optimization are used to craft malicious training sets that evade standard detection methods [148]. Federated learning [148] further increases these risks, as attackers may inject poisoned weights and biases into the central server or transmit them from multiple decentralized nodes. Adversarial data manipulation [149] also extends to inference, where crafted inputs [150] bypass detection. Both black-box [151] and white-box [152] methods reveal model vulnerabilities, while backdoor triggers [153] quietly activate malicious behaviors in trained models.

Adversaries also refine publicly available research or software tools [154] [144] to craft advanced strategies. A 2021 demonstration by Kaspersky [155] showed how adversaries could adapt local feature extraction methods to bypass cloud-based malware detection. Similarly, poisoned datasets publication in open repositories [60] represent another vector, as attackers alter data or labels [19] before organizations

download them. Real-world incidents, such as the VirusTotal poisoning [156] and Microsoft’s Tay chatbot compromise [157], reveal how subtle manipulations degrade detection engines or produce offensive behaviors.

In conclusion, data integrity attacks undermine the foundations of ML systems by compromising the data used during training and inference. Poisoning, adversarial manipulation, and the publication of tampered datasets can lead to degraded performance, hidden vulnerabilities, and a loss of trust [87]. Addressing these risks requires thorough validation, secure data pipelines, and a clear understanding of how adversaries embed malicious content. Building on these concerns, the next section examines model integrity-level attacks, which target the internal structure and behavior of ML models.

### E. Model Integrity Attacks

Model integrity attacks target the training, architecture, and deployment of ML systems, exploiting weaknesses that can undermine accuracy, fairness, and reliability. They manifest through various techniques, including data poisoning, adversarial input data manipulation, and backdoor embedding. Each of them are intended to corrupt model’s behavior or outputs. These attacks can degrade decision-making, expose private data, and spread misinformation. Robust controls are necessary to safeguard models, especially as collaborative development and open-source resources become more prevalent.

Backdoor attacks [98] are a prominent technique, embedding malicious triggers through poisoned data [149], or direct payloads [15]. The model appears normal until specific inputs activate hidden functionalities. Eroding model’s integrity [83] causes end-users to lose trust and confidence in the system over time. It coerces the victim organization to lose money as well as time to repair the model, a job that otherwise ought to have been automated. A 2023 demonstration [158] showed how uploading a tainted large language model to HuggingFace posed supply chain risks. Similarly, direct payload injection [109] modifies the model itself, causing unexpected responses when triggered. In some computer vision tasks [148], adversaries incorporate visual markers to mislead classification, while generative or reinforcement learning models have been tricked into producing harmful outputs.

Attacks that exploit full ML model access [159] give adversaries white-box knowledge of architectures, and/or parameters. Armed with this knowledge, they inject tailored backdoors or craft adversarial inputs that undermine the system. A study [88] showed how “neural payload injection” affected real-world Android apps from Google Play, revealing 54 vulnerable applications used for cash recognition, parental control, face authentication, and financial services. These backdoors could remain hidden if deployment environments lack proper safeguards. Publishing poisoned models [62] is another strategy, where compromised models with dormant triggers are placed in repositories. If unsuspecting developers integrate these models, the triggers stay inactive until malicious conditions arise. Meanwhile, ML supply chain compromises [160] occur when shared models lack proper provenance, letting attackers introduce biases or malware.

In summary, model integrity level attacks arise at every phase of the ML lifecycle, from training and distribution to deployment. Outdated or deprecated models [15] also pose risks when they disregard security updates, opening the door for manipulations like ROME. Addressing these threats calls for secure development practices, rigorous model testing, and stronger provenance controls. The next area of concern involves deployment runtime-level attacks, which target the environments where models operate.

### F. Deployment Runtime Attacks

Deployment runtime attacks involve targeting ML models during their operational phase, where adversaries exploit real-time data exchanges and user interactions [71]. They often submit carefully crafted adversarial inputs [149] to inference APIs [161], revealing weaknesses like overfitting or susceptibility to subtle modifications that appear harmless to human observers. Such methods enable resource exhaustion [15] [84] or misclassification, ultimately undermining the reliability of MLOps ecosystem. Over time, these attacks erode trust [83] in automated decision-making, emphasizing the need for continuous monitoring and proactive defenses. In addition, stealthy manipulations can bypass detection [79], introducing long-term vulnerabilities that are challenging to pinpoint once integrated into normal workflows.

Adversaries also exfiltrate data through inference APIs [162], enabling them to extract model behavior or replicate proprietary logic [108] [163] [164]. These incursions threaten data confidentiality [15], compromise intellectual property [86], and jeopardize system stability [15]. Attacks on LLMs frequently involve prompt injection [113] [15], embedding adversarial directives [109] [110] that manipulate outputs or enable unauthorized actions. Universal adversarial triggers [165] further show how concise token sequences can override model contexts, generating hateful or misleading content when triggered. In some cases, employees have inadvertently shared confidential information with LLMs [166], illustrating the risk of data leakage [167] [15]. Additionally, self-replication attacks [168], such as those used in the Morris II worm [114], reveal how malicious prompts can propagate across interconnected systems, amplifying risk.

Meta prompt extraction [69] [15] enables attackers to uncover hidden system instructions and internal logic, potentially exposing proprietary knowledge or subverting intended behaviors. Adversaries often validate these exploits using verify attacks [169] in low-risk environments before targeting production systems. This approach allows them to fine-tune attack methods while minimizing detection. ML-enabled products and services [170] are also frequent targets. For example, malware detection models can be bypassed [79] by appending specific token sequences or leveraging cloned translation models to craft adversarial queries, thereby degrading performance and compromising intellectual property [86].

A notable example highlighting such risk is the attack on public machine translation services [171], where researchers from UC Berkeley replicated near state-of-the-art translation service using black-box APIs from Google Translate and

Systran Translate. By training a surrogate model through repeated queries, they demonstrated how intellectual property can be functionally extracted from these systems. They then deployed adversarial inputs crafted for the surrogate model [165], successfully transferring word flips, vulgar content, and omitted sentences to the actual production systems. This case study underscores the wide range of runtime threats, from operational disruption to the theft and misuse of proprietary algorithms. Mitigating such attacks requires a layered defense strategy, incorporating continuous input validation, robust monitoring, and user education. As organizations increasingly adopt shared models and foundational architectures, a single compromise may propagate across the ecosystem. The next section examines network exploitation-level attacks, which target the underlying communication infrastructure that supports machine learning operations.

### G. Network Exploitation Attacks

Network exploitation attacks target the interconnected infrastructure of machine learning systems, aiming to disrupt availability [81], compromise data [172], or degrade pipelines [172]. By exploiting communication pathways, exposed endpoints, and interfaces, adversaries can manipulate ML workflows or undermine trust in the MLOps ecosystem [173] [174]. The interconnected nature of modern ML architectures amplifies these attacks, highlighting the need for strong security controls to safeguard shared resources.

Public-facing applications [174] are a common entry point. In the ShadowRay incident [5], attackers exploited Ray’s Jobs API, which lacked proper authentication, to execute arbitrary commands. Denial of ML service [81] [15] is another technique, where adversaries flood ML systems with malicious requests to deplete computational resources. Similarly, spamming systems with chaff data [82] [172] overwhelms detection mechanisms by introducing noise, causing overfitting or wasted analyst attention. These methods erode reliability and increase costs if not addressed.

Data exfiltration [173] via network either through API or web-application further endangers MLOps ecosystem, as attackers target datasets, models, and proprietary algorithms. In one case, malicious scripts embedded in Google Colab notebooks [6] granted unauthorized access to ML artifacts stored in cloud environments, exposing intellectual property and weakening competitiveness. Unbounded consumption [15] also poses risks, particularly for large language models, by enabling unauthorized use of inference capabilities and imposing resource burdens. Meanwhile, adversaries exploit vulnerable protocols [175] or launch reflection and amplification attacks [176] to manipulate data in transit or magnify malicious traffic.

Endpoints in distributed ML architectures represent another attack vector as it can result in unauthorized access, allowing adversaries to manipulate or disrupt ML workflows. Attackers may hijack web sessions [177], forge credentials [178], or spoof DHCP configurations [179] to gain control over sensitive components. Network traffic manipulation [172] can corrupt data pipelines by injecting poisoned or malformed data, causing model degradation or delays. The broader implications

of network-level exploitation attacks extend to the overall security posture of ML systems. As adversaries continue to develop sophisticated techniques, the attack surface within MLOps ecosystem expands. The dynamic nature of these systems, coupled with their reliance on distributed architectures, necessitates a proactive approach to security.

Network exploitation represents just one facet of the broader adversarial landscape within the MLOps ecosystem. Other categories such as exploitation, access abuse, pipeline infrastructure, data integrity, model integrity, and deployment runtime attacks underscore the wide array of vulnerabilities present throughout the lifecycle. Exploitation attacks leverage social engineering or deceptive prompts to influence user behavior, while access abuse involves the misuse of valid credentials or privileged accounts to infiltrate ML systems. Pipeline infrastructure attacks compromise the workflows that enable ML development and deployment, whereas data integrity attacks corrupt the datasets that models rely on. Model integrity attacks tamper with the internal logic or behavior of ML models, and deployment runtime attacks disrupt operational environments post-deployment. Having examined these diverse attack vectors, the following section introduces mitigation strategies aimed at strengthening ML systems against such attacks, promoting secure, reliable, and resilient operations across the evolving MLOps landscape.

## V. SECURITY MITIGATION STRATEGIES

The previous section outlined the adversarial landscape and identified key vectors that attackers may exploit. This section presents corresponding mitigation strategies, drawing from the MITRE ATLAS framework, MITRE ATT&CK, and additional relevant literature. Each subsection aligns mitigation approaches with the attack categories discussed earlier, offering a structured defense strategy for addressing threats introduced in Section IV. Since most attacks begin with a preparatory phase before targeting users or systems directly, it is essential to consider these early-stage activities in the mitigation process. Table I summarizes preparatory tactics, including Reconnaissance, Resource Development, Discovery, Collection, and Impact, along with associated mitigation techniques. The table also includes hyperlinks to guide readers to the relevant subsections. We do not provide a separate subsection for these preparatory mitigations, as many are already covered within the techniques described in this section or represent opportunities for future research.

We begin by addressing mitigations for user-level attacks. Fig. 6 summarizes key mitigation strategies, grouped into two main categories: Exploitation and Access Abuse. Each category highlights practical techniques that end users can adopt to reduce exposure to specific threats. The next subsection expands on the mitigation strategies related to exploitation-based attacks introduced in Fig. 6.

### A. Exploitation Mitigations

Mitigating exploitation-level attacks requires coordinated technical controls, user training, and proactive safeguards.

TABLE I

AN OVERVIEW OF THE MAPPING OF THE TACTICS AND TECHNIQUES DRAWN FROM THE SECTION III AND MITIGATIONS DRAWN FROM THE MITRE ATLAS. SOME TECHNIQUES REQUIRE FURTHER RESEARCH ANALYSIS, AND SOME ARE REFERRED TO IN SECTION V.

| Tactic               | Techniques   | Mitigations  |
|----------------------|--|--|
| Reconnaissance       | Search for Victim’s Publicly Available Research Materials, Search Victim-Owned Websites, Search Application Repositories | Limit Public Release of Information [180]  |
|                      | Active Scanning  | Potential Research Gap   |
| Resource Development | Discussed in Section IV  | Discussed in Section V   |
| Discovery            | Discover ML Model Ontology, Discover ML Model Family   | Restrict Number of ML Model Queries [181]  |
|                      | Discover ML Model Family   | Use Ensemble Methods [182]   |
|                      | LLM Meta Prompt Extraction   | Discussed in Section V-F   |
|                      | Discover LLM Hallucinations, Discover AI Model Outputs   | Potential Research Gap   |
| Collection           | ML Artifact Collection   | Encrypt Sensitive Information [183]  |
|                      | Data from Information Repositories, Data from Local System   | Potential Research Gap   |
| Impact               | Evade ML Model   | Discussed in Section V-F   |
|                      | Denial of ML Service   | Discussed in Section V-G   |
|                      | Erode ML Model Integrity   | Discussed in Section V-E   |
|                      | Spamming ML System with Chaff Data   | Discussed in Section V-G   |
|                      | Cost Harvesting  | Restrict Number of ML Model Queries [181]  |
|                      | External Harms: ML Intellectual Property Theft   | Control Access to ML Models and Data at Rest [184],<br>Encrypt Sensitive Information [183] |
|                      | External Harms, Erode Dataset Integrity  | Potential Research Gap   |

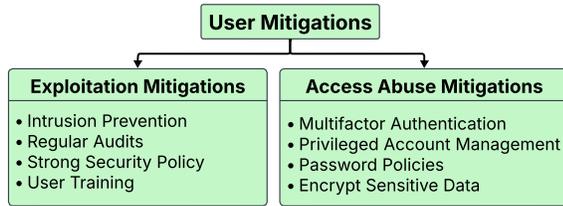


Fig. 6. Overview of key mitigation strategies against user-level attacks.

Organizations can reduce phishing exposure by deploying antivirus and antimalware tools [185], conducting regular audits and scans [186], and blocking suspicious traffic using Intrusion Prevention Systems [187]. Restricting web-based content [188] and securing software configurations [189] further limit attack surfaces. Ongoing user education [190] [191] enhances awareness of social engineering, while developers are trained in secure coding and ML-specific vulnerabilities. To counter user execution attacks, enforcing restricted library loading [192], using code signing [193], verifying ML artifacts [194], applying vulnerability scanning [195], and maintaining an AI Bill of Materials [196] strengthens supply chain integrity.

Mitigating hallucinated content involves refining inputs, model architectures, and outputs [197]. Retrieval-Augmented Generation (RAG), Knowledge Retrieval, and the Decompose and Query framework validate outputs using external sources [197]. Real-time Verification and Rectification (EVER) [197] verifies model outputs during generation. Post-generation strategies such as Retrofit Attribution using Research and Revision (RARR) and High Entropy Word Spotting and Replacement clean outputs after they are generated. Self-refinement approaches such as Prompting GPT-3 To Be Reliable, ChatProtect, and Self-Reflection leverage the model’s internal reasoning [197]. Structured Comparative (SC) reasoning, Mind’s Mirror, Chain-of-Verification (CoVe), and Chain of Natural Language Inference (CoNLI) introduce step-by-step validation workflows. Prompt engineering methods including UPRISE and SynTra optimize instructions to reduce hallucinations. Decoding strategies such as Context-Aware Decoding (CAD), Decoding by Contrasting Layers (DoLa),

and Inference-Time Intervention (ITI) influence generation to favor factual content. Knowledge graph-based tools like RHO and FLEEK [197] embed structured knowledge into responses. Additional methods include loss function tuning (THAM, Loss Weighting), knowledge injection, teacher-student learning, and Hallucination-Augmented Recitations (HAR) reinforce model grounding. Factuality can also be improved through fine-tuning and Refusal-Aware Instruction Tuning (R-Tuning) that help models recognize knowledge boundaries, while Think While Effectively Articulating Knowledge (TWEAK) treats outputs as hypotheses requiring verification [197].

Privilege escalation risks can be mitigated through strict security policies, the use of multifactor authentication, and endpoint lockdown. Biometric techniques such as keystroke dynamics provide an additional layer of behavioral authentication, helping to detect abnormal activity early. Maintaining documented security policies, keeping tools up to date, and controlling all system connections further strengthen resilience. Applying least privilege principles, separating duties, and incorporating behavioral analytics helps limit abuse. Some organizations also rely on temporal signals and non-Markovian models to identify anomalous behavior over time [198]. Together, these strategies address common exploitation scenarios, including phishing, hallucination, privilege escalation, and user execution. The next subsection outlines mitigation techniques aimed at access abuse-level attacks.

### B. Access Abuse Mitigations

Mitigating access abuse attacks involves controlling and monitoring user accounts to prevent unauthorized access and reduce potential damage from compromised credentials. This includes strict access policies, authentication methods, and continuous oversight of privileged and regular user activities. Account use policies [199] and Active Directory configurations [200] help define access guidelines and enforce secure authentication and authorization practices. Multi-factor authentication [201] further strengthens mitigation strategies, while password policies [202] and privileged account man-

agement [203] reduce vulnerabilities. User account management [204] supports effective access control, and application developer guidance [205] further reduces the risk of security gaps. Regular user training [190] raises awareness and deters common social engineering techniques.

Protecting against unsecured credentials adds another layer of defense. Active Directory configuration [200] and regular audits [186] identify possible weaknesses, and encrypting sensitive information [183] alongside network traffic filtering [206] can limit exposure. Restricting resource access over networks [207] and applying strict operating system configurations [208] reduce the attack surface. Enforcing password policies [202] and managing privileged accounts [203] are also key, along with restricting file and directory permissions [209]. Routine software updates [210] and user training [190] strengthen this strategy by ensuring that systems stay current and personnel remain vigilant. These approaches address the risk of valid accounts being misused or credentials left unprotected. The next subsection will examine mitigation techniques focused on system-level attacks as shown in Fig. 7, beginning with strategies for securing pipeline infrastructure.

### C. Pipeline Infrastructure Mitigations

Securing the MLOps pipeline, from data acquisition to model deployment, involves protecting infrastructure, data, and model artifacts throughout their lifecycle. Key mitigation strategies in this area are shown in Fig. 7. Pre-compromise measures [211] aim to secure environments before adversaries gain a foothold, since no existing defenses prevent them from acquiring infrastructure for attacks, such as purchasing servers or domains, or using free resources like expired domains or virtual servers to evade attribution. Ongoing research focuses on detecting the misuse of open-source tools, expired domains, and AI-generated scripts for malicious purposes.

Mitigating GPU hardware compromise requires layered security measures [212]. Techniques include hardware-based memory encryption (e.g., Intel’s Total Memory Encryption and AMD’s Secure Memory Encryption) and secure key management using Hardware Security Modules and Trusted Platform Modules [212]. Memory integrity is maintained through hash trees and Message Authentication Codes. Mitigations strategies against side-channel include cache partitioning and memory access randomization. Other key mitigations involve speculative execution control, physical isolation of enclaves, and machine learning-based mitigations against evolving attacks. Firmware integrity can be maintained through cryptographic signature verification, tamper-resistant physical protections, and real-time monitoring of anomalies in voltage, frequency, and temperature [212]. Redundancy and error correction mechanisms, including EDC and ECC, enhance fault tolerance, while shielding, randomized delays, and electromagnetic filtering help defend against injection attacks. Additional safeguards include automated security policy enforcement, memory access randomization, and secure enclave execution [212].

Buffer overflow mitigations span both hardware and software mechanisms [213]. Hardware-level protections such as

No-Execute (NX) bits, Address Space Layout Randomization (ASLR), Position Independent Executables (PIE), and read-only relocations prevent code injection. Techniques like StackGuard, SmashGuard, and secure return address stacks protect return addresses through either hardware or software enforcement [213]. Trusted hardware-based defenses, including eXecute Only Memory (XOM) and Secure Bit, are further reinforced by hybrid solutions such as HSDDefender. Compiler-level strategies include ProPolice for stack layout protection, StackShield for return address duplication, Return Address Defender for maintaining protected copies, and Address Sanitizer for detecting memory errors. Automatic Fortification additionally enforces buffer bounds checking, contributing to a multi-layered defense against overflow vulnerabilities [213].

To prevent ML software compromise, ensemble methods [182] improve resistance to adversarial inputs, while code signing [193] ensures software integrity across the supply chain. For command and scripting interpreter attacks, multiple layers of mitigation are necessary. Antivirus and antimalware tools [185], system audits [186], behavior prevention on endpoints [214], and execution prevention [215] all contribute to this effort. Disabling unused features [216], limiting software installation [217], and managing privileged accounts [204] reduce potential attack vectors. Restricting web-based content [188] limits additional entry points. Physical environment access risks can be reduced with multi-modal sensors [218], which combine various data sources to detect unauthorized access and tampering.

For persistence mitigation, systems like Cyber Persistence Detector (CPD) [219] track potential persistence mechanisms by linking setup and execution events using pseudo-dependency edges. Expert-guided edges and detection rules help identify suspicious activities. The system correlates related kill chain events, generates concise attack graphs, and adjusts weighting factors to reflect varying attack contexts. An alert budget system prioritizes key events for investigation, focusing on behavioral patterns rather than easily altered indicators. This approach improves detection and mitigation of persistence attacks in enterprise environments. In summary, these techniques address a range of attacks on ML pipelines, spanning infrastructure, hardware, software, and physical access. The following subsection outlines mitigation strategies aimed at preserving data integrity and limiting attacker control within ML environments.

### D. Data Integrity Mitigations

Protecting data integrity in ML systems demands verifying dataset authenticity, guarding against poisoning, and maintaining detailed provenance. Diverse mitigation strategies focus on adversarial ML attacks by removing redundant or sensitive data through deduplication and sanitization [220]. Techniques such as differential privacy inject noise to protect sensitive information [221], while encryption-based approaches secure gradients and model outputs [220]. Randomized smoothing, data augmentation, and model ensembling improve resilience to input perturbations [222], while adding friendly noise makes adversarial example generation difficult [223]. Reinforcement

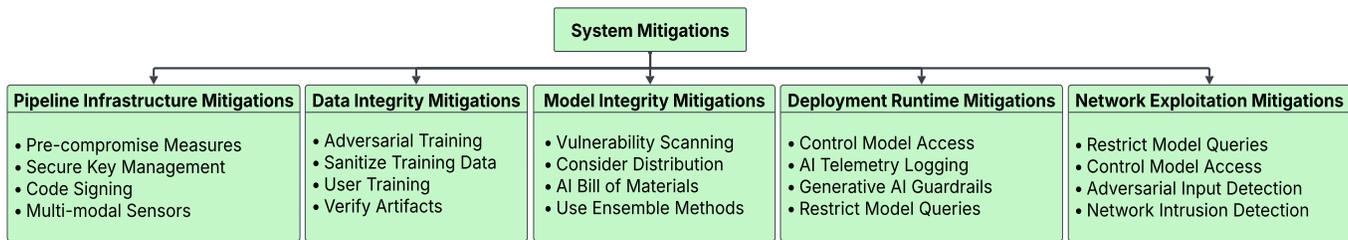


Fig. 7. Overview of key mitigation strategies against system-level attacks.

learning from human feedback helps align models with acceptable behavior [220]. Data anonymization and end-to-end encryption safeguard personal information [220], whereas proof-of-work puzzles can slow extraction attempts [224]. Profiling natural examples and setting energy consumption thresholds limit the impact of attacks aimed at maximizing resource usage [225]. Regular patching addresses discovered vulnerabilities [226], and user education clarifies attack risks [220].

To preserve datasets integrity, researchers use multiple methods to identify and deter risks during data collection. Researchers can adopt mixed collection methods, comparing closed distributions with verified contact lists to open distributions for broader reach. Timing marketing efforts effectively, employing survey-specific fraud tests, and piloting those tests before rollout can improve response quality [227]. Paradata (such as submission times, geolocation, and device details) reveals anomalies, while baseline samples from trusted participants provide benchmarks for unusual behaviors [227]. Carefully designed open-ended questions, including those tailored as “fraud catchers” [227], help highlight suspicious responses. Incentive structures may include physical or location-based rewards to discourage mass participation. Careful review of free-text responses, along with checks for duplicate IPs, shared contacts, unlikely answer patterns, and closely timed submissions, can help identify inconsistencies [227]. Automated scoring solutions like MinFraud Score detect suspicious entries, and continuous monitoring of these indicators is advised to account for adaptive attacks. Verification emails or phone calls can confirm legitimacy, and in contexts of high online fraud, paper or mail surveys may serve as an alternative [227].

Verifying ML artifacts [194] through cryptographic checksums mitigates poisoned dataset risks, while an AI Bill of Materials [196] traces data sources and modifications. This approach provides a record of all artifacts, capturing each step in the data lifecycle for supply chain risk management. Limiting the release of model artifacts and controlling data access [184] restrict unauthorized tampering, while sanitizing training data [228] identifies and removes compromised inputs. Maintaining AI dataset provenance [229] ensures visibility into data origins and changes, offering a method to compare current datasets with known trusted versions.

Adversarial training, gradient masking, and input preprocessing counter many manipulations by either strengthening the model against attempted exploits or filtering unusual inputs before inference [230]. Defensive distillation smooths model boundaries, making minor perturbations less effective. Randomized smoothing or enforcing Lipschitz continuity can yield

certified robustness within defined thresholds [230]. Detection tools draw on statistical anomaly techniques or specialized auxiliary models to flag suspicious samples. As attackers develop new approaches, explainable mitigation strategies allow teams to interpret and respond to emerging patterns, and automated solutions help update defense strategies without manual intervention [230]. Collaboration between ML and cybersecurity experts promotes research that merges attack intelligence with model engineering, supporting standards and guidelines for organizations deploying AI solutions [230]. These approaches safeguard data quality and authenticity while documenting all data source. The next subsection will explore techniques focused on defending model integrity level attacks.

### E. Model Integrity Mitigations

Protecting machine learning models from tampering and unauthorized alterations requires several measures throughout the model’s lifecycle. One approach involves scanning model artifacts for potential vulnerabilities [195], which is particularly useful for detecting exploits in file formats like pickle files. Controlling access to model registries and restricting production model visibility to approved users also reduces the risk of manipulation [184]. Sanitizing training data [228] helps detect, remove, or filter poisoned inputs, including explicit or unwanted content that may compromise model performance. This approach can be especially relevant for active learning environments, where newly ingested data should be filtered. Regular validation of models [231] can further reveal backdoors, adversarial bias, or data poisoning through tests for concept drift or shifts in training data. In addition, maintaining a complete record of dataset sources and changes [229] contributes to transparency and traceability across different stages of model development.

Organizations can also manage how models are deployed and shared to mitigate exposure. Serving models in the cloud instead of on edge devices can limit the level of access adversaries have [232], and computing features in the cloud can reduce gray-box attacks that stem from knowledge of preprocessing techniques. An AI Bill of Materials [196] enhances supply chain awareness by listing all artifacts that went into building the AI system, which allows faster responses to reported vulnerabilities. Using ensemble methods [182] can increase robustness against adversarial inputs, while code signing [193] helps maintain software integrity. These strategies can be combined to mitigate scenarios where a model is compromised, including cases in which poisoned models are published or attackers gain extensive control over the model.

In summary, these measures are designed to protect models from unauthorized modification, compromised distribution, or the inclusion of malicious components. The next subsection focuses on mitigation techniques for deployment runtime attacks, where the model is actively served in production.

#### F. Deployment Runtime Mitigations

Organizations can address deployment runtime attacks by combining access controls, monitoring, and secure software development practices. Enforcing robust authentication for production APIs, requiring user identity verification, and restricting ML model or data access in production [233] limits malicious use. Restricting query controls further reduces adversaries opportunities to overload systems or extract information. Monitoring model queries and behaviors in real time with collecting AI telemetry data [234] enables organizations to detect suspicious patterns, uncover potential exfiltration attempts, and maintain an audit trail for incident response.

ML models can also be shielded by hardening techniques [235], such as adversarial training, which helps the model learn from known adversarial inputs. Using an ensemble [182] of diverse models likewise prevents single points of failure because an attack’s effectiveness against one architecture may fail against another. Using multiple sensors [218] and pre-processing inputs [236] can reduce the impact of physical attacks or perturbations introduced by adversaries. Additionally, detecting anomalous or malicious queries early [237] by evaluating request patterns against known legitimate behaviors helps block adversarial attempts at the entry point.

Guarding large language model (LLM) deployments requires systematic guardrails [238] placed between user interactions and the generated output. These can take the form of filters, rule-based logic, or AI-driven classifiers to confirm that prompts and responses meet security requirements. Generative AI Guidelines [239], often appended to user prompts or embedded into system instructions, define the acceptable scope of model behavior and clarify safety parameters. Aligning Generative AI models [240] through supervised fine-tuning or reinforcement learning introduces additional checks during training, ensuring that alignment with organizational policy remains intact despite model updates.

Multiple LLM-specific attacks benefit from these measures. Prompt injection attacks and data leakage risks can be addressed by restricting unauthorized content, logging model inputs [234], and blocking suspicious queries [238]. Meta prompt extraction, which involves revealing hidden instructions or chain-of-thought data, can also be combated through generative AI guidelines [239] and alignment strategies [240]. LLM plugins that access external systems require close oversight to prevent compromise; consistent use of guardrails keeps the plugin’s interactions within pre-approved boundaries [238].

For LLM prompt self-replication, including LLM Tagging [241] can identify the origin of responses, while Delimiting Data [241] and random sequence enclosure [241] make it harder for prompts to replicate across interactions. The Sandwich Defense layers user instructions around prior responses, and Instruction Defense [241] clarifies that models should

never alter user inputs. Combining these with Marking [241], which distinguishes user prompts from agent outputs, can significantly reduce cross-model injection attempts. However, reliance on advanced models can introduce new risks, especially if these models themselves become compromised.

Organizations should also guard against verify attacks by restricting offline copies of ML models and watching for suspicious API usage [220]. Regularly patching known vulnerabilities and educating developers and users about potential exploitation tactics can help in early attack detection. Combining this with thorough logging and frequent audits closes many attack surfaces that attackers could exploit [226].

#### G. Network Exploitation Mitigations

To address attacks at the network level, organizations can secure communication channels, limit unauthorized access, and protect network infrastructure that supports AI operations [13]. Public-facing applications benefit from strict input validation, sanitization, and least privilege access [242]. Regular software updates, use of Web Application Firewalls [174], and TLS-based secure communication [174] help prevent exploitation attempts. Security headers, rate limiting, thorough logging, and consistent monitoring further strengthen defenses. Practices such as application isolation [243], exploit protection [244], network segmentation [245], and privileged account management [203] add layers of control and visibility. Frequent vulnerability scans [246] also reduce exposure.

Organizations can mitigate denial of ML service by limiting the volume and rate of queries [181] and filtering adversarial inputs [237], while restricting queries counters spamming with chaff data [181]. To reduce exfiltration risks, access to internal model registries and production models should be controlled [184]. Network intrusion prevention systems [187] can block suspicious traffic, adding protection against data leakage.

By implementing these measures, organizations improve security against attacks on AI infrastructure such as data exfiltration and denial-of-service. Regular software updates [210], vulnerability scans [246], and privileged access controls [203] provide additional layers of defense, while continuous monitoring and timely response remain essential for long-term protection. In this section, we discussed mitigation strategies that address user exploitation, access abuse, pipeline infrastructure exploitations, data and model integrity concerns, deployment runtime risks, and network exploitation. Together, these methods safeguard various stages of MLOps from user compromise and data gathering to model training and deployment. Table II provides a concise overview of the taxonomy categories, including hyperlinks to corresponding attacks and mitigation strategies, and indicates their relevance to specific phases of the MLOps ecosystem. This tabular representation enables readers to clearly understand the relationship between the identified attack surfaces, the applicable mitigation strategies, and the stages of the MLOps lifecycle they impact. Despite these efforts, mitigation strategies remain limited in scope. In the next section, we propose research challenges and recommendations that can further help strengthen MLOps security across these interconnected domains.

TABLE II  
AN OVERVIEW OF THE MAPPING OF THE ATTACKS DRAWN FROM THE SECTION IV AND AND MITIGATIONS DRAWN FROM THE SECTION V TO THE MLOPS PHASES AS IDENTIFIED IN THE FIG. 1.

| Taxonomy                | Attacks | Mitigations | MLOps Phases |                 |                   |                   |                  |            |
|-------------------------|---------|-------------|--------------|-----------------|-------------------|-------------------|------------------|------------|
|                         |         |             | Admin Setup  | Data Collection | Model Development | Approval Workflow | Model Deployment | Monitoring |
| Exploitation            | [IV-A]  | [V-A]       | ✓            | ✓               |                   | ✓                 |                  |            |
| Access Abuse            | [IV-B]  | [V-B]       | ✓            |                 |                   | ✓                 |                  | ✓          |
| Pipeline Infrastructure | [IV-C]  | [V-C]       | ✓            | ✓               | ✓                 | ✓                 | ✓                | ✓          |
| Data Integrity          | [IV-D]  | [V-D]       |              | ✓               | ✓                 |                   |                  |            |
| Model Integrity         | [IV-E]  | [V-E]       |              |                 | ✓                 |                   |                  |            |
| Deployment Runtime      | [IV-F]  | [V-F]       |              |                 |                   |                   | ✓                |            |
| Network Exploitation    | [IV-G]  | [V-G]       | ✓            | ✓               |                   |                   | ✓                | ✓          |

## VI. RESEARCH CHALLENGES AND RECOMMENDATION

Although the previous section outlined a range of mitigation strategies to defend the MLOps ecosystem, these measures often provide only partial protection against increasingly sophisticated, AI-driven adversaries. As attacks evolve rapidly, organizations must address not only technical vulnerabilities but also legal, ethical, and operational complexities. The following discussion highlights key concerns and proposes future research directions, with the goal of strengthening the security of the MLOps ecosystem while preserving essential principles of transparency, privacy, and compliance.

### A. The Rise of AI-Driven Social Engineering

Attackers now exploit generative AI to craft highly convincing phishing emails, clone voices for vishing schemes, and produce deepfake videos that entice employees into revealing sensitive credentials [247]. Once inside an organization, these adversaries can move laterally and compromise vital MLOps components. Although zero trust models and ongoing user education are vital, they frequently fall behind AI’s rapid advancement. Further research should refine deepfake detection techniques, enhance digital authentication protocols, and implement comprehensive incident monitoring to counter the increasingly realistic nature of AI-enabled social engineering.

### B. Malicious Repositories and LLM Hallucinations

LLMs occasionally generate misleading “hallucinations”, which can enable attackers to insert malicious packages into trusted online repositories. If organizations inadvertently integrate these compromised resources, they risk exposing intellectual property, disrupting workflows, or introducing harmful payloads. Threat actors may also upload cloned or tampered models to platforms such as Hugging Face, deceiving unsuspecting practitioners into deploying compromised versions [248]. Future research should explore automated systems for detecting suspicious content, robust versioning protocols to ensure code integrity, and community-driven validation frameworks supported by the broader MLOps ecosystem.

### C. Misuse of AI Tools by Advanced Persistent Threats (APTs)

Advanced Persistent Threats (APTs) groups sponsored by nation state actors increasingly rely on freely available powerful AI tools, such as Google’s Gemini, to identify vulnerabilities, develop sophisticated malware, and evade standard security measures [249]. Google DeepMind’s Gemma team has raised concerns about potential misuse of AI in upcoming

releases [250]. Researchers should examine subtle ways AI may be weaponized, establish precise ethical and legal standards, and propose robust prevention strategies. Effectively defending MLOps ecosystem against increasingly sophisticated adversaries will require detecting and preventing abuse of emerging AI technologies.

### D. Benchmarking MLOps Defenses

Many firms rely on laboratory-style tests that do not accurately reflect real-world attacks. As a result, MLOps ecosystem remain vulnerable when deployed in practical settings. Future research should focus on benchmarking defenses against adaptive and diverse attack patterns [251]. By simulating adversarial scenarios, such as targeted phishing attempts and sophisticated AI-powered attacks, MLOps teams can evaluate how well their security measures perform under pressure. Additionally, continuous research should establish standardized criteria for comparing defensive mechanisms. This approach will support ongoing improvement and ensure better alignment with real operational risks.

### E. Red-Teaming for Robust Security

Red team exercises serve as an effective approach for uncovering hidden vulnerabilities, validating existing defenses, and strengthening overall system resilience. Drawing on initiatives such as MITRE’s red team program, organizations can more accurately simulate real-world attack scenarios and identify exploit paths that may remain undetected during conventional testing. Future research should explore frameworks that support recurring, collaborative red-teaming efforts across sectors, ensuring that shared insights contribute to a more secure global MLOps ecosystem. By embracing these rigorous assessments, security teams can enhance threat detection, improve incident response, and cultivate a proactive security posture.

### F. Strengthening MLOps Robustness

A robust MLOps strategy adopts security as a default posture and embeds high-assurance practices throughout the ML lifecycle. Regular testing using established adversarial simulation tools and tailored threat scenarios can help uncover and remediate vulnerabilities before they lead to critical failures. Further investigation is needed to identify optimal strategies for streamlining configuration management, continuous patching, and resilience assessments across the MLOps ecosystem. This proactive approach enables organizations to anticipate, disrupt, and adapt to evolving adversarial threats.

### G. Open Source vs. Closed Source

An important question in securing MLOps is how to balance openness with security when comparing closed-source and open-source approaches. Publicly accessible code fosters collaborative innovation but may also provide adversaries with insights into system architecture. Conversely, limited disclosure can hinder attackers but may reduce transparency and erode user trust. Further research should explore how community-driven governance, selective disclosure, and cryptographic verification can preserve the benefits of open development while reducing reconnaissance risks. By striking this balance, organizations can maintain stakeholder trust and safeguard operational integrity.

### H. Cultivating Cyber-Hygiene Among MLOps Practitioners

Developers, DevOps engineers, and data scientists often work in silos, unintentionally introducing security risks in MLOps environments. While password encryption and secure networking are crucial, practitioners also need training in secure coding [190], proactive vulnerability management, and regular data backups. Further research is needed to assess the effectiveness of training methods, incentive structures, and policies in fostering a security-oriented culture at both organizational and societal levels.

### I. Balancing Transparency with Trust

Organizations often publish white papers and technical reports to showcase their commitment to ethical and responsible AI. However, excessive disclosure of MLOps infrastructure can provide adversaries with a blueprint for exploitation. Researchers must develop frameworks that promote transparency in fairness criteria and privacy protections without revealing sensitive operational details. Future research should focus on defining indicators of responsible AI that offer public and regulatory assurance without compromising system security. A well-balanced transparency model can help organizations build trust while safeguarding confidential information.

### J. The Challenge of Copyright, Privacy, and Compliance

Generative AI models may unintentionally expose copyrighted or sensitive personal data, leading to significant legal consequences and reputational damage. When models continuously train on real or user-provided data, issuing disclaimers provides only a minimal level of protection [252]. Scholars should develop comprehensive solutions that include frequent model audits, automated data validation, and tightly controlled input and output streams. These measures can help AI advancements comply with legal requirements while safeguarding user privacy and intellectual property. Further research should enhance auditing processes, establish measurable compliance benchmarks, and design robust fail-safes to minimize liability in high-risk AI applications.

Protecting MLOps ecosystem requires a holistic strategy that combines technological solutions with organizational best practices. By embracing rigorous benchmarking, conducting thorough red-team assessments, and establishing frameworks that prioritize security from the onset, organizations can better defend their MLOps ecosystems against emerging attacks.

## VII. CONCLUSION

In recent years, MLOps has fundamentally reshaped how machine learning models are developed, deployed, and maintained. It offers substantial advantages, including enhanced scalability, reproducibility, rapid deployment, and improved team collaboration. However, the acceleration of development to deployment cycles, combined with increased automation and highly interconnected workflows, significantly expands the attack surface. This heightens the risk of vulnerabilities and security breaches. As machine learning becomes increasingly embedded in critical systems, these trends underscore the urgent need for cybersecurity frameworks that can evolve alongside technological advancements.

This study highlights that while MLOps delivers critical operational benefits, it also introduces security challenges that demand focused attention. Leveraging the MITRE ATLAS framework for AI-specific threats, we present a structured taxonomy of vulnerabilities spanning the entire MLOps lifecycle, from data collection to deployment. Insights from red-team exercises and real-world incidents illustrate how adversaries exploit both user-level and system-level weaknesses. To address these threats, we propose targeted mitigation strategies informed by established frameworks and best practices. These recommendations offer actionable guidance for improving MLOps security and identifying priorities for future research.

As MLOps adoption continues to grow, organizations must embrace a proactive security posture by embedding robust practices early in their machine learning workflows. By implementing the frameworks and recommendations outlined in this work, they can more effectively safeguard models, maintain stakeholder confidence, and ensure that machine learning continues to support secure, ethical, and responsible innovation. In doing so, they establish a resilient foundation for long-term security, sustained innovation, and adaptive growth in a rapidly evolving threat landscape.

## REFERENCES

- [1] D. Sculley, G. Holt, D. Golovin, E. Davydov, T. Phillips, D. Ebner, V. Chaudhary, M. Young, J.-F. Crespo, and D. Dennison, "Hidden technical debt in Machine learning systems," in *Proceedings of the 29th International Conference on Neural Information Processing Systems - Volume 2*, NIPS'15, (Cambridge, MA, USA), MIT Press, 2015.
- [2] T. Mucci and C. Stryker, "What is MLOps? — IBM." <https://www.ibm.com/think/topics/mlops>, Apr 2024. Accessed: 2025-04-10.
- [3] National Institute of Standards and Technology, "Development Operations (DevOps) - Glossary." [https://csrc.nist.gov/glossary/term/development\\_operations](https://csrc.nist.gov/glossary/term/development_operations). Accessed: 2025-04-10.
- [4] J. Stone, R. Patel, F. Ghiasi, S. Mittal, and S. Rahimi, "Navigating MLOps: Insights into Maturity, Lifecycle, Tools, and Careers." <https://arxiv.org/abs/2503.15577>, 2025.
- [5] A. Lumelsky, G. Kaplan, and G. Elbaz, "ShadowRay: First Known Attack Campaign Targeting AI Workloads Actively Exploited In The Wild." <https://www.oligo.security/blog/shadowray-attack-ai-workloads-actively-exploited-in-the-wild>, Mar 2024. Accessed: 2025-01-11.
- [6] 4n7m4n, "Careful Who You Colab With: Abusing Google Colaboratory." <https://antman1p-30185.medium.com/careful-who-you-colab-with-fa8001f933e7>, Jul 2022. Accessed: 2025-01-26.
- [7] D.-A. team, "DeepSeek-V3 Technical Report." <https://arxiv.org/abs/2412.19437>, 2024.

- [8] M. Gault, "OpenAI Claims DeepSeek Plagiarized Its Plagiarism Machine." <https://gizmodo.com/openai-claims-deepseek-plagiarized-its-plagiarism-machine-2000556339>, Jan 2025. Accessed: 2025-02-23.
- [9] Z. Whittaker, "Security lapse exposed Clearview AI source code," *TechCrunch*, Apr 2020. Accessed: 2025-02-02.
- [10] S. Cyber, "Cylance, I Kill You!." <https://skylightcyber.com/2019/07/18/cylance-i-kill-you/>, Jul 2019. Accessed: 2025-02-01.
- [11] B. Yu, J. Pan, J. Hu, A. Nascimento, and M. De Cock, "Character Level based Detection of DGA Domain Names," in *2018 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8, 2018.
- [12] K. Baker, "Open Source Intelligence (OSINT)." <https://www.crowdstrike.com/en-us/cybersecurity-101/threat-intelligence/open-source-intelligence-osint/>, Jan 2025. Accessed: 2025-02-01.
- [13] "MITRE ATLAS™." <https://atlas.mitre.org/>. Accessed: 2024-07-27.
- [14] MITRE Corporation, "MITRE ATT&CK Framework." <https://attack.mitre.org/>, 2024. Accessed: 2024-12-08.
- [15] OWASP, "OWASP Top 10 for LLM Applications 2025." <https://owasp.org/www-project-top-10-for-large-language-model-applications/assets/PDF/OWASP-Top-10-for-LLMs-v2025.pdf>, Nov 2024. Accessed: 2025-01-11.
- [16] V. Tolpegin, S. Truex, M. E. Gursoy, and L. Liu, "Data poisoning attacks against federated learning systems," in *Computer security—ESORICs 2020: 25th European symposium on research in computer security, ESORICs 2020, guildford, UK, Sep 14–18, 2020, proceedings, part i 25*, pp. 480–501, Springer, 2020.
- [17] M. Testi, M. Ballabio, E. Frontoni, G. Iannello, S. Moccia, P. Soda, and G. Vessio, "MLOps: a taxonomy and a methodology," *IEEE Access*, vol. 10, pp. 63606–63618, 2022.
- [18] S. Mäkinen, H. Skogström, E. Laaksonen, and T. Mikkonen, "Who needs MLOps: What data scientists seek to accomplish and how can MLOps help?," in *2021 IEEE/ACM 1st Workshop on AI Engineering-Software Engineering for AI (WAIN)*, pp. 109–112, IEEE, 2021.
- [19] F. A. Yerlikaya and Şerif Bahtiyar, "Data poisoning attacks against machine learning algorithms," *Expert Systems with Applications*, vol. 208, p. 118101, 2022.
- [20] D. Kreuzberger, N. Kühl, and S. Hirschl, "Machine learning operations (mlops): Overview, definition, and architecture," *IEEE Access*, vol. 11, pp. 31866–31879, 2023.
- [21] X. Zhang and J. Jaskolka, "Conceptualizing the secure machine learning operations (secmlops) paradigm," in *2022 IEEE 22nd International Conference on Software Quality, Reliability and Security (QRS)*, pp. 127–138, IEEE, 2022.
- [22] I. Kumara, R. Arts, D. Nucci, W. Heuvel, and D. Tamburri, "Requirements and reference architecture for mlops:insights from industry," 11 2022.
- [23] Amazon, Inc., "AWS Identity and Access Management." <https://aws.amazon.com/iam/>. Accessed: 2024-7-14.
- [24] M. M. John, H. H. Olsson, and J. Bosch, "Towards MLOps: A Framework and Maturity Model," pp. 334–341, Institute of Electrical and Electronics Engineers Inc., 9 2021.
- [25] A. Lima, L. Monteiro, and A. P. Furtado, "MLOps: Practices, Maturity Models, Roles, Tools, and Challenges - A Systematic Literature Review," vol. 1, Science and Technology Publications, Lda, 2022.
- [26] S. Garg, P. Pundir, G. Rathee, P. Gupta, S. Garg, and S. Ahlawat, "On Continuous Integration / Continuous Delivery for Automated Deployment of Machine Learning Models using MLOps," in *2021 IEEE Fourth International Conference on Artificial Intelligence and Knowledge Engineering (AIKE)*, pp. 25–28, 2021.
- [27] J. Kazmierczak, K. Salama, V. Huerta, and S. K. J. Bahadur, "MLOps: Continuous delivery and automation pipelines in machine learning," Aug 2024.
- [28] A. Ali, R. Pinciroli, F. Yan, and E. Smirmi, "BATCH: Machine Learning Inference Serving on Serverless Platforms with Adaptive Batching," in *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*, pp. 1–15, 2020.
- [29] "NVIDIA: MLPerf AI Benchmarks." <https://www.nvidia.com/en-us/data-center/resources/mlperf-benchmarks/>. Accessed: 2024-07-27.
- [30] V. J. e. a. Reddi, "MLPerf Inference Benchmark," in *2020 ACM/IEEE 47th Annual International Symposium on Computer Architecture (ISCA)*, pp. 446–459, 2020.
- [31] N. Ferry, A. Rossini, F. Chauvel, B. Morin, and A. Solberg, "Towards Model-Driven Provisioning, Deployment, Monitoring, and Adaptation of Multi-cloud Systems," in *2013 IEEE Sixth International Conference on Cloud Computing*, pp. 887–894, 2013.
- [32] MITRE ATT&CK, "Enterprise Tactics." <https://attack.mitre.org/tactics/enterprise/>, 2024. Accessed: 2025-02-01.
- [33] MITRE ATT&CK, "Enterprise Techniques." <https://attack.mitre.org/techniques/enterprise/>, 2025. Accessed: 2025-02-01.
- [34] MITRE ATLAS, "Reconnaissance." <https://atlas.mitre.org/tactics/AML.TA0002>, 2024. Accessed: 2025-02-01.
- [35] MITRE ATLAS, "Search for Victim's Publicly Available Research Materials: Journals and Conference Proceedings." <https://atlas.mitre.org/techniques/AML.T0000.000>, May 2021. Accessed: 2025-02-01.
- [36] MITRE ATLAS, "Search for Victim's Publicly Available Research Materials: Pre-Print Repositories." <https://atlas.mitre.org/techniques/AML.T0000.001>, May 2021. Accessed: 2025-02-01.
- [37] MITRE ATLAS, "Search for Victim's Publicly Available Research Materials: Technical Blogs." <https://atlas.mitre.org/techniques/AML.T0000.002>, May 2021. Accessed: 2025-02-01.
- [38] MITRE ATLAS, "Search for Publicly Available Adversarial Vulnerability Analysis." <https://atlas.mitre.org/techniques/AML.T0001>, Mar 2024. Accessed: 2025-02-01.
- [39] MITRE ATLAS, "Search Victim-Owned Websites." <https://atlas.mitre.org/techniques/AML.T0003>, May 2021. Accessed: 2025-02-01.
- [40] MITRE ATLAS, "Search Application Repositories." <https://atlas.mitre.org/techniques/AML.T0004>, May 2021. Accessed: 2025-02-01.
- [41] MITRE ATLAS, "Active Scanning." <https://atlas.mitre.org/techniques/AML.T0006>, Jan 2023. Accessed: 2025-02-01.
- [42] Y. Miao, C. Chen, L. Pan, Q.-L. Han, J. Zhang, and Y. Xiang, "Machine learning-based cyber attacks targeting on controlled information: A survey," *ACM Computing Surveys (CSUR)*, vol. 54, pp. 1–36, 2021.
- [43] MITRE ATLAS, "Resource Development." <https://atlas.mitre.org/tactics/AML.TA0003>, Jan 2023. Accessed: 2025-02-01.
- [44] MITRE ATLAS, "Acquire Public ML Artifacts." <https://atlas.mitre.org/techniques/AML.T0002>, May 2021. Accessed: 2025-02-01.
- [45] MITRE ATLAS, "Acquire Public ML Artifacts: Datasets." <https://atlas.mitre.org/techniques/AML.T0002.000>, 2021. Accessed: 2025-02-01.
- [46] MITRE ATLAS, "Acquire Public ML Artifacts: Models." <https://atlas.mitre.org/techniques/AML.T0002.001>, 2021. Accessed: 2025-02-01.
- [47] N. P. et al., "Technical Report on the CleverHans v2.1.0 Adversarial Examples Library," *arXiv preprint arXiv:1610.00768*, 2018.
- [48] I. Research, *Adversarial Robustness Toolbox*, Jul 2020.
- [49] J. Rauber, R. Zimmermann, M. Bethge, and W. Brendel, "Foolbox Native: Fast adversarial attacks to benchmark the robustness of machine learning models in PyTorch, TensorFlow, and JAX," *Journal of Open Source Software*, vol. 5, no. 53, p. 2607, 2020.
- [50] MITRE ATLAS, "Obtain Capabilities." <https://atlas.mitre.org/techniques/AML.T0016>, 2023. Accessed: 2025-02-01.
- [51] B. E. Granger and F. Pérez, "Jupyter: Thinking and Storytelling With Code and Data," *Computing in Science & Engineering*, pp. 7–14, 2021.
- [52] MITRE ATLAS, "Develop Capabilities." <https://atlas.mitre.org/techniques/AML.T0017>, Oct 2023. Accessed: 2025-02-01.
- [53] MITRE ATLAS, "Acquire Infrastructure." <https://atlas.mitre.org/techniques/AML.T0008>, May 2021. Accessed: 2025-02-01.
- [54] MITRE ATLAS, "Acquire Infrastructure: Physical Countermeasures." <https://atlas.mitre.org/techniques/AML.T0008.003>, Oct 2024. Accessed: 2025-02-01.
- [55] E. Bisong, *Google Colaboratory*. Berkeley, CA: Apress, 2019.
- [56] F. P. Miller, A. F. Vandome, and J. McBrewster, *Amazon Web Services*. Alpha Press, 2010.
- [57] M. Corporation, "Microsoft azure." <https://azure.microsoft.com/>. Accessed: 2025-02-02.
- [58] MITRE ATLAS, "Acquire Infrastructure: ML Development Workspaces." <https://atlas.mitre.org/techniques/AML.T0008.000>, May 2021. Accessed: 2025-02-01.
- [59] MITRE ATLAS, "Acquire Infrastructure: Domains." <https://atlas.mitre.org/techniques/AML.T0008.002>, Oct 2024. Accessed: 2025-02-01.
- [60] MITRE ATLAS, "Publish Poisoned Datasets." <https://atlas.mitre.org/techniques/AML.T0019>, May 2021. Accessed: 2025-02-01.
- [61] MITRE ATLAS, "Poison Training Data." <https://atlas.mitre.org/techniques/AML.T0020>, May 2021. Accessed: 2025-02-01.
- [62] MITRE ATLAS, "Publish Poisoned Models." <https://atlas.mitre.org/techniques/AML.T0058>, Oct 2024. Accessed: 2025-02-01.
- [63] MITRE ATLAS, "Publish Hallucinated Entities." <https://atlas.mitre.org/techniques/AML.T0060>, Oct 2024. Accessed: 2025-02-01.
- [64] MITRE ATLAS, "Establish Accounts." <https://atlas.mitre.org/techniques/AML.T0021>, Jan 2023. Accessed: 2025-02-01.
- [65] MITRE ATLAS, "Discovery." <https://atlas.mitre.org/tactics/AML.TA0008>, Jan 2023. Accessed: 2025-02-01.
- [66] MITRE ATLAS, "Discover ML Artifacts." <https://atlas.mitre.org/techniques/AML.T0007>, May 2021. Accessed: 2025-02-01.
- [67] MITRE ATLAS, "Discover ML Model Ontology." <https://atlas.mitre.org/techniques/AML.T0013>, May 2021. Accessed: 2025-02-01.

- [68] MITRE ATLAS, “Discover ML Model Family.” <https://atlas.mitre.org/techniques/AML.T0014>, May 2021. Accessed: 2025-02-01.
- [69] MITRE ATLAS, “LLM Meta Prompt Extraction.” <https://atlas.mitre.org/techniques/AML.T0056>, Oct 2023. Accessed: 2025-02-01.
- [70] MITRE ATLAS, “Discover LLM Hallucinations.” <https://atlas.mitre.org/techniques/AML.T0062>, Oct 2024. Accessed: 2025-02-01.
- [71] MITRE ATLAS, “Discover AI Model Outputs.” <https://atlas.mitre.org/techniques/AML.T0063>, Oct 2024. Accessed: 2025-02-01.
- [72] MITRE ATLAS, “Face Identification System Evasion via Physical Countermeasures.” <https://atlas.mitre.org/studies/AML.CS0012>, Jan 2020. Accessed: 2025-02-01.
- [73] MITRE ATLAS, “Collection.” <https://atlas.mitre.org/tactics/AML.TA0009>, Jan 2023. Accessed: 2025-02-01.
- [74] MITRE ATLAS, “ML Artifact Collection.” <https://atlas.mitre.org/techniques/AML.T0035>, May 2021. Accessed: 2025-02-01.
- [75] MITRE ATLAS, “Data from Information Repositories.” <https://atlas.mitre.org/techniques/AML.T0036>, Jan 2023. Accessed: 2025-02-01.
- [76] MITRE ATLAS, “Data from Local System.” <https://atlas.mitre.org/techniques/AML.T0037>, Jan 2023. Accessed: 2025-02-01.
- [77] MITRE ATLAS, “Microsoft Azure Service Disruption.” <https://atlas.mitre.org/studies/AML.CS0010>, 2020. Accessed: 2025-02-01.
- [78] MITRE ATLAS, “Impact.” <https://atlas.mitre.org/tactics/AML.TA0011>, Jan 2023. Accessed: 2025-02-01.
- [79] MITRE ATLAS, “Evade ML Model.” <https://atlas.mitre.org/techniques/AML.T0015>, Oct 2022. Accessed: 2025-02-01.
- [80] P. Olson, “Faces Are the Next Target for Fraudsters,” *The Wall Street Journal*, Jul 2021. Accessed: 2025-02-02.
- [81] MITRE ATLAS, “Denial of ML Service.” <https://atlas.mitre.org/techniques/AML.T0029>, May 2021. Accessed: 2025-02-01.
- [82] MITRE ATLAS, “Spamming ML System with Chaff Data.” <https://atlas.mitre.org/techniques/AML.T0046>, 2021. Accessed: 2025-02-01.
- [83] MITRE ATLAS, “Erode ML Model Integrity.” <https://atlas.mitre.org/techniques/AML.T0031>, May 2021. Accessed: 2025-02-01.
- [84] MITRE ATLAS, “Cost Harvesting.” <https://atlas.mitre.org/techniques/AML.T0034>, May 2021. Accessed: 2025-02-01.
- [85] MITRE ATLAS, “External Harms.” <https://atlas.mitre.org/techniques/AML.T0048>, Oct 2023. Accessed: 2025-02-01.
- [86] MITRE ATLAS, “External Harms: ML Intellectual Property Theft.” <https://atlas.mitre.org/techniques/AML.T0048.004>, Oct 2023. Accessed: 2025-02-01.
- [87] MITRE ATLAS, “Erode Dataset Integrity.” <https://atlas.mitre.org/techniques/AML.T0059>, Oct 2024. Accessed: 2025-02-01.
- [88] Y. Li, J. Hua, H. Wang, C. Chen, and Y. Liu, “DeepPayload: Black-box Backdoor Attack on Deep Learning Models through Neural Payload Injection.” <https://arxiv.org/abs/2101.06896>, 2021.
- [89] F. McKenna, “The Many Jobs and Wigs of Eric Jaklitch’s Fraud Scheme,” *Frank on Fraud*, Jan 2022. Accessed: 2025-02-02.
- [90] MITRE ATLAS, “Phishing.” <https://atlas.mitre.org/techniques/AML.T0052>, Oct 2023. Accessed: 2025-02-01.
- [91] MITRE ATT&CK, “Phishing.” <https://attack.mitre.org/techniques/T1566/>, Oct 2024. Accessed: 2025-02-02.
- [92] Computronix, “AI Phishing Attacks: How Cybercriminals Use Automation to Target You.” <https://computronixusa.com/ai-phishing-attacks-cybercriminals-leveraging-automation/>, 2025. Accessed: 2025-01-26.
- [93] MITRE ATLAS, “Phishing: Spearphishing via Social Engineering LLM.” <https://atlas.mitre.org/techniques/AML.T0052.000>, Oct 2023. Accessed: 2025-02-01.
- [94] MITRE ATLAS, “Valid Accounts.” <https://atlas.mitre.org/techniques/AML.T0012>, Jan 2023. Accessed: 2025-02-01.
- [95] MITRE ATT&CK, “Firmware.” <https://attack.mitre.org/datasources/DS0001/>, Mar 2022. Accessed: 2025-02-02.
- [96] N. Carlini, M. Jagielski, C. A. Choquette-Choo, D. Paleka, W. Pearce, H. Anderson, A. Terzis, K. Thomas, and F. Tramèr, “Poisoning Web-Scale Training Datasets is Practical,” in *2024 IEEE Symposium on Security and Privacy (SP)*, pp. 407–425, 2024.
- [97] N. Pitropakis, E. Panaousis, T. Giannetos, E. Anastasiadis, and G. Loukas, “A taxonomy and survey of attacks against machine learning,” *Computer Science Review*, vol. 34, p. 100199, 2019.
- [98] MITRE ATLAS, “Backdoor ML Model.” <https://atlas.mitre.org/techniques/AML.T0018>, May 2021. Accessed: 2025-02-01.
- [99] MITRE ATLAS, “Create Proxy ML Model.” <https://atlas.mitre.org/techniques/AML.T0005>, May 2021. Accessed: 2025-02-01.
- [100] MITRE ATT&CK, “Spoof Reporting Message.” <https://attack.mitre.org/techniques/T0856/>, Oct 2023. Accessed: 2025-02-02.
- [101] MITRE ATT&CK, “Browser Session Hijacking.” <https://attack.mitre.org/techniques/T1185/>, Feb 2022. Accessed: 2025-02-02.
- [102] G. Sebastian and P. Kolluru, “Rethinking the Weakest Link in the Cybersecurity Chain,” *ISACA Journal*, 2021. Accessed: 2025-02-02.
- [103] MITRE ATLAS, “Privilege Escalation.” <https://atlas.mitre.org/tactics/AML.TA0012>, Oct 2023. Accessed: 2025-02-01.
- [104] MITRE ATLAS, “User Execution.” <https://atlas.mitre.org/techniques/AML.T0011>, Jan 2023. Accessed: 2025-02-01.
- [105] MITRE ATLAS, “Unsecured Credentials.” <https://atlas.mitre.org/techniques/AML.T0010>, May 2021. Accessed: 2025-02-01.
- [106] MITRE ATLAS, “User Execution: Unsafe ML Artifacts.” <https://atlas.mitre.org/techniques/AML.T0011.000>, 2021. Accessed: 2025-02-01.
- [107] MITRE ATLAS, “User Execution: Malicious Package.” <https://atlas.mitre.org/techniques/AML.T0011.001>, 2024. Accessed: 2025-02-01.
- [108] MITRE ATLAS, “Exfiltration via ML Inference API: Infer Training Data Membership.” <https://atlas.mitre.org/techniques/AML.T0024.000>, May 2021. Accessed: 2025-02-01.
- [109] MITRE ATLAS, “LLM Prompt Injection: Direct.” <https://atlas.mitre.org/techniques/AML.T0051.000>, Oct 2023. Accessed: 2025-02-01.
- [110] MITRE ATLAS, “LLM Prompt Injection: Indirect.” <https://atlas.mitre.org/techniques/AML.T0051.001>, Oct 2023. Accessed: 2025-02-01.
- [111] MITRE ATLAS, “LLM Plugin Compromise.” <https://atlas.mitre.org/techniques/AML.T0053>, Oct 2023. Accessed: 2025-02-01.
- [112] MITRE ATLAS, “LLM Jailbreak.” <https://atlas.mitre.org/techniques/AML.T0054>, Oct 2023. Accessed: 2025-02-01.
- [113] MITRE ATLAS, “LLM Prompt Injection.” <https://atlas.mitre.org/techniques/AML.T0051>, Oct 2023. Accessed: 2025-02-01.
- [114] S. Zakeer, “Self-replicating Morris II worm targets ai email assistants.” <https://securityintelligence.com/posts/morris-ii-self-replicating-malware-genai-email-assistants/>, Jun 2024.
- [115] Wunderwuzzi, “ChatGPT Plugins: Data Exfiltration via Images & Cross Plugin Request Forgery.” <https://embracethered.com/blog/posts/2023/chatgpt-webpilot-data-exfil-via-markdown-injection/>, May 2023. Accessed: 2025-01-26.
- [116] MITRE ATT&CK, “Credential Access.” <https://attack.mitre.org/tactics/TA0006/>, Jul 2019. Accessed: 2025-03-02.
- [117] MITRE ATT&CK, “Valid Accounts.” <https://attack.mitre.org/techniques/T1078/>, Oct 2024. Accessed: 2025-02-02.
- [118] MITRE ATT&CK, “Valid Accounts: Default Accounts.” <https://attack.mitre.org/techniques/T1078/001/>, Oct 2024. Accessed: 2025-02-02.
- [119] MITRE ATT&CK, “Valid Accounts: Cloud Accounts.” <https://attack.mitre.org/techniques/T1078/004/>, Oct 2024. Accessed: 2025-02-02.
- [120] MITRE ATT&CK, “Valid Accounts: Local Accounts.” <https://attack.mitre.org/techniques/T1078/003/>, Oct 2024. Accessed: 2025-02-02.
- [121] MITRE ATLAS, “Unsecured Credentials.” <https://atlas.mitre.org/techniques/AML.T0055>, Apr 2024. Accessed: 2025-02-01.
- [122] MITRE ATT&CK, “Unsecured Credentials.” <https://attack.mitre.org/techniques/T1552/>, Oct 2024. Accessed: 2025-02-02.
- [123] MITRE ATT&CK, “Unsecured Credentials: Cloud Instance Metadata API.” <https://attack.mitre.org/techniques/T1552/005/>, Oct 2024. Accessed: 2025-02-02.
- [124] E. A. Brewer, “Kubernetes and the path to cloud native,” in *Proceedings of the Sixth ACM Symposium on Cloud Computing, SoCC '15*, (New York, NY, USA), p. 167, Association for Computing Machinery, 2015.
- [125] D. Merkel, “Docker: lightweight Linux containers for consistent development and deployment,” *Linux J.*, vol. 2014, Mar. 2014.
- [126] MITRE ATT&CK, “Unsecured Credentials: Container API.” <https://attack.mitre.org/techniques/T1552/007/>, 2024. Accessed: 2025-02-02.
- [127] MITRE ATT&CK, “Unsecured Credentials: Credentials in Registry.” <https://attack.mitre.org/techniques/T1552/002/>, 2024. Accessed: 2025-02-02.
- [128] MITRE ATT&CK, “Unsecured Credentials: Bash History.” <https://attack.mitre.org/techniques/T1552/003/>, 2024. Accessed: 2025-02-02.
- [129] MITRE ATT&CK, “Unsecured Credentials: Private Keys.” <https://attack.mitre.org/techniques/T1552/004/>, 2024. Accessed: 2025-02-02.
- [130] Slack, “Slack.” <https://slack.com/>. Accessed: 2025-02-02.
- [131] Atlassian, “Trello — Project Management Software — trello.com.” <https://trello.com/>, 2024. [Accessed 2024-12-16].
- [132] M. Corporation, “Microsoft teams.” <https://www.microsoft.com/en-us/microsoft-teams/group-chat-software>. Accessed: 2025-02-02.
- [133] D. Hendrycks, C. Burns, S. Kadavath, A. Arora, S. Basart, E. Tang, D. Song, and J. Steinhardt, “Measuring Mathematical Problem Solving With the MATH Dataset.” <https://arxiv.org/abs/2103.03874>, 2021.
- [134] L. Cazorla, C. Alcaraz, and J. Lopez, “Cyber Stealth Attacks in Critical Information Infrastructures,” *IEEE Systems Journal*, 2018.
- [135] MITRE ATLAS, “Acquire Infrastructure: Consumer Hardware.” <https://atlas.mitre.org/techniques/AML.T0008.001>, 2021. Accessed: 2025-02-01.

- [136] MITRE ATLAS, "ML Supply Chain Compromise: Hardware." <https://atlas.mitre.org/techniques/AML.T0010.000>, 2021. Accessed: 2025-02-01.
- [137] MITRE ATLAS, "Physical Environment Access." <https://atlas.mitre.org/techniques/AML.T0041>, May 2021. Accessed: 2025-02-01.
- [138] MITRE ATLAS, "Command and Scripting Interpreter." <https://atlas.mitre.org/techniques/AML.T0050>, Oct 2023. Accessed: 2025-02-01.
- [139] G. van Rossum, "Python tutorial," Tech. Rep. CS-R9526, Centrum voor Wiskunde en Informatica (CWI), Amsterdam, May 1995.
- [140] Microsoft Corporation, *PowerShell*. Accessed: 2025-02-05.
- [141] MITRE, "Common Attack Pattern Enumeration and Classification (CAPEC)." <https://capec.mitre.org>. [Online; accessed 2025-01-11].
- [142] P. Kocher, J. Horn, A. Fogh, D. Genkin, D. Gruss, W. Haas, M. Hamburg, M. Lipp, S. Mangard, T. Prescher, M. Schwarz, and Y. Yarom, "Spectre Attacks: Exploiting Speculative Execution," in *2019 IEEE Symposium on Security and Privacy (SP)*, pp. 1–19, 2019.
- [143] MITRE ATLAS, "ML Supply Chain Compromise: ML Software." <https://atlas.mitre.org/techniques/AML.T0010.001>, 2021. Accessed: 2025-02-01.
- [144] MITRE ATLAS, "Obtain Capabilities: Software Tools." <https://atlas.mitre.org/techniques/AML.T0016.001>, 2023. Accessed: 2025-02-01.
- [145] A. Sharma, "PyTorch discloses malicious dependency chain compromise over holidays," Jan 2023. Accessed: 2025-01-26.
- [146] MITRE ATLAS, "Persistence." <https://atlas.mitre.org/tactics/AML.TA0006>, Jan 2023. Accessed: 2025-02-01.
- [147] S. Chen, J. Xu, E. C. Sezer, P. Gauriar, and R. K. Iyer, "Non-control-data attacks are realistic threats," in *USENIX security symposium*, vol. 5, p. 146, 2005.
- [148] M. Goldblum, D. Tsipras, C. Xie, X. Chen, A. Schwarzschild, D. Song, A. Madry, B. Li, and T. Goldstein, "Dataset Security for Machine Learning: Data Poisoning, Backdoor Attacks, and Defenses," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [149] MITRE ATLAS, "Craft Adversarial Data." <https://atlas.mitre.org/techniques/AML.T0043>, May 2021. Accessed: 2025-02-01.
- [150] A. Kurakin, I. Goodfellow, and S. Bengio, "Adversarial examples in the physical world!" <https://arxiv.org/abs/1607.02533>, 2017.
- [151] MITRE ATLAS, "Craft Adversarial Data: Black-Box Optimization." <https://atlas.mitre.org/techniques/AML.T0043.001>, May 2021. Accessed: 2025-02-01.
- [152] MITRE ATLAS, "Craft Adversarial Data: White-Box Optimization." <https://atlas.mitre.org/techniques/AML.T0043.000>, May 2021. Accessed: 2025-02-01.
- [153] MITRE ATLAS, "Craft Adversarial Data: Insert Backdoor Trigger." <https://atlas.mitre.org/techniques/AML.T0043.004>, May 2021. Accessed: 2025-02-01.
- [154] MITRE ATLAS, "Develop Capabilities: Adversarial ML Attacks." <https://atlas.mitre.org/techniques/AML.T0017.000>, 2023. Accessed: 2025-02-01.
- [155] A. Antonov and A. Kogtenkov, "How to confuse antimalware neural networks. Adversarial attacks and protection," Jun 2021. Accessed: 2025-01-26.
- [156] MITRE ATLAS, "VirusTotal Poisoning." <https://atlas.mitre.org/studies/AML.CS0002>, 2020. Accessed: 2025-02-01.
- [157] O. Schwartz, "In 2016, Microsoft's Racist Chatbot Revealed the Dangers of Online Conversation," Nov 2019. Accessed: 2025-01-26.
- [158] D. Huynh and J. Hardouin, "PoisonGPT: How We Hid a Lobotomized LLM on Hugging Face to Spread Fake News," *Mithril Security Blog*, Jul 2023. Accessed: 2025-02-05.
- [159] MITRE ATLAS, "Full ML Model Access." <https://atlas.mitre.org/techniques/AML.T0044>, May 2021. Accessed: 2025-02-01.
- [160] MITRE ATLAS, "ML Supply Chain Compromise: Model." <https://atlas.mitre.org/techniques/AML.T0010.003>, May 2021. Accessed: 2025-02-01.
- [161] MITRE ATLAS, "AI Model Inference API Access." <https://atlas.mitre.org/techniques/AML.T0040>, Oct 2024. Accessed: 2025-02-01.
- [162] MITRE ATLAS, "Exfiltration via ML Inference API." <https://atlas.mitre.org/techniques/AML.T0024>, May 2021. Accessed: 2025-02-01.
- [163] MITRE ATLAS, "Exfiltration via ML Inference API: Invert ML Model." <https://atlas.mitre.org/techniques/AML.T0024.001>, May 2021. Accessed: 2025-02-01.
- [164] MITRE ATLAS, "Exfiltration via ML Inference API: Extract ML Model." <https://atlas.mitre.org/techniques/AML.T0024.002>, May 2021. Accessed: 2025-02-01.
- [165] E. Wallace, S. Feng, N. Kandpal, M. Gardner, and S. Singh, "Universal Adversarial Triggers for Attacking and Analyzing NLP." <https://arxiv.org/abs/1908.07125>, 2021.
- [166] V. Petkauskas, "Lessons learned from ChatGPT's Samsung leak." <https://cybernews.com/security/chatgpt-samsung-leak-explained-lessons/>, May 2023.
- [167] MITRE ATLAS, "LLM Data Leakage." <https://atlas.mitre.org/techniques/AML.T0057>, Oct 2023. Accessed: 2025-02-01.
- [168] MITRE ATLAS, "LLM Prompt Self-Replication." <https://atlas.mitre.org/techniques/AML.T0061>, Oct 2024. Accessed: 2025-02-01.
- [169] MITRE ATLAS, "Verify Attack." <https://atlas.mitre.org/techniques/AML.T0042>, May 2021. Accessed: 2025-02-01.
- [170] MITRE ATLAS, "ML-Enabled Product or Service." <https://atlas.mitre.org/techniques/AML.T0047>, May 2021. Accessed: 2025-02-01.
- [171] E. Wallace, M. Stern, and D. Song, "Imitation Attacks and Defenses for Black-box Machine Translation Systems," 2021.
- [172] B. I. Rubinstein, B. Nelson, L. Huang, A. D. Joseph, S.-h. Lau, S. Rao, N. Taft, and J. D. Tygar, "ANTIDOTE: Understanding and defending against poisoning of anomaly detectors," in *Proceedings of the 9th ACM SIGCOMM Conference on Internet Measurement, IMC '09*, (New York, NY, USA), p. 1–14, Association for Computing Machinery, 2009.
- [173] MITRE ATLAS, "Exfiltration via Cyber Means." <https://atlas.mitre.org/techniques/AML.T0025>, May 2021. Accessed: 2025-02-01.
- [174] MITRE ATLAS, "Exploit Public-Facing Application." <https://atlas.mitre.org/techniques/AML.T0049>, Feb 2023. Accessed: 2025-02-01.
- [175] MITRE ATT&CK, "Application Layer Protocol: Web Protocols." <https://attack.mitre.org/techniques/T1107/001/>, 2024. Accessed: 2025-02-02.
- [176] MITRE ATT&CK, "Network Denial of Service: Reflection Amplification." <https://attack.mitre.org/techniques/T1498/002/>, Oct 2024. Accessed: 2025-02-02.
- [177] MITRE ATT&CK, "Browser Session Hijacking." <https://attack.mitre.org/techniques/T1185/>, Feb 2022. Accessed: 2025-02-02.
- [178] MITRE ATT&CK, "Forge Web Credentials." <https://attack.mitre.org/techniques/T1606/>, Oct 2024. Accessed: 2025-02-02.
- [179] MITRE ATT&CK, "Adversary-in-the-Middle: DHCP Spoofing." <https://attack.mitre.org/techniques/T1557/003/>, 2024. Accessed: 2025-02-02.
- [180] MITRE ATLAS, "Limit Public Release of Information." <https://atlas.mitre.org/mitigations/AML.M0000>, Oct 2024. Accessed: 2025-02-01.
- [181] MITRE ATLAS, "Restrict Number of ML Model Queries." <https://atlas.mitre.org/mitigations/AML.M0004>, Oct 2023. Accessed: 2025-02-01.
- [182] MITRE ATLAS, "Use Ensemble Methods." <https://atlas.mitre.org/mitigations/AML.M0023>, Oct 2023. Accessed: 2025-02-01.
- [183] MITRE ATT&CK, "Encrypt Sensitive Information." <https://attack.mitre.org/mitigations/M1041/>, Jun 2019. Accessed: 2025-02-02.
- [184] MITRE ATLAS, "Control Access to ML Models and Data at Rest." <https://atlas.mitre.org/mitigations/AML.M0005>, 2023. Accessed: 2025-02-01.
- [185] MITRE ATT&CK, "Antivirus/Antimalware." <https://attack.mitre.org/mitigations/M1049/>, Mar 2020. Accessed: 2025-02-02.
- [186] MITRE ATT&CK, "Audit." <https://attack.mitre.org/mitigations/M1047/>, Oct 2024. Accessed: 2025-02-02.
- [187] MITRE ATT&CK, "Network Intrusion Prevention." <https://attack.mitre.org/mitigations/M1031/>, Oct 2024. Accessed: 2025-02-02.
- [188] MITRE ATT&CK, "Restrict Web-Based Content." <https://attack.mitre.org/mitigations/M1021/>, Jun 2019. Accessed: 2025-02-02.
- [189] MITRE ATT&CK, "Software Configuration." <https://attack.mitre.org/mitigations/M1054/>, Dec 2023. Accessed: 2025-02-02.
- [190] MITRE ATLAS, "User Training." <https://atlas.mitre.org/mitigations/AML.M0018>, Oct 2023. Accessed: 2025-02-01.
- [191] MITRE ATT&CK, "User Training." <https://attack.mitre.org/mitigations/M1017/>, Oct 2024. Accessed: 2025-02-02.
- [192] MITRE ATLAS, "Restrict Library Loading." <https://atlas.mitre.org/mitigations/AML.M0011>, Oct 2023. Accessed: 2025-02-01.
- [193] MITRE ATLAS, "Code Signing." <https://atlas.mitre.org/mitigations/AML.M0013>, Oct 2023. Accessed: 2025-02-01.
- [194] MITRE ATLAS, "Verify ML Artifacts." <https://atlas.mitre.org/mitigations/AML.M0014>, Oct 2023. Accessed: 2025-02-01.
- [195] MITRE ATLAS, "Vulnerability Scanning." <https://atlas.mitre.org/mitigations/AML.M0016>, Jan 2024. Accessed: 2025-02-01.
- [196] MITRE ATLAS, "AI Bill of Materials." <https://atlas.mitre.org/mitigations/AML.M0023>, Oct 2024. Accessed: 2025-02-01.
- [197] S. Tonmoy, S. Zaman, V. Jain, A. Rani, V. Rawte, A. Chadha, and A. Das, "A comprehensive survey of hallucination mitigation techniques in large language models," *arXiv:2401.01313*, 2024.
- [198] M. Mehmood, R. Amin, M. M. A. Muslam, J. Xie, and H. Aldabbas, "Privilege escalation attack detection and mitigation in cloud using machine learning," *IEEE Access*, vol. 11, pp. 46561–46576, 2023.
- [199] MITRE ATT&CK, "Account Use Policies." <https://attack.mitre.org/mitigations/M1036/>, Oct 2022. Accessed: 2025-02-02.

- [200] MITRE ATT&CK, "Active Directory Configuration." <https://attack.mitre.org/mitigations/M1015/>, Oct 2024. Accessed: 2025-02-02.
- [201] MITRE ATT&CK, "Multi-factor Authentication." <https://attack.mitre.org/mitigations/M1032/>, Oct 2022. Accessed: 2025-02-02.
- [202] MITRE ATT&CK, "Password Policies." <https://attack.mitre.org/mitigations/M1027/>, Oct 2022. Accessed: 2025-02-02.
- [203] MITRE ATT&CK, "Privileged Account Management." <https://attack.mitre.org/mitigations/M1026/>, Oct 2024. Accessed: 2025-02-02.
- [204] MITRE ATT&CK, "User Account Management." <https://attack.mitre.org/mitigations/M1018/>, May 2020. Accessed: 2025-02-02.
- [205] MITRE ATT&CK, "Application Developer Guidance." <https://attack.mitre.org/mitigations/M1013/>, Sep 2023. Accessed: 2025-02-02.
- [206] MITRE ATT&CK, "Filter Network Traffic." <https://attack.mitre.org/mitigations/M1037/>, Oct 2024.
- [207] MITRE ATT&CK, "Limit Access to Resource Over Network." <https://attack.mitre.org/mitigations/M1035/>, 2020. Accessed: 2025-02-02.
- [208] MITRE ATT&CK, "Operating System Configuration." <https://attack.mitre.org/mitigations/M1028/>, Mar 2023. Accessed: 2025-02-02.
- [209] MITRE ATT&CK, "Restrict File and Directory Permissions." <https://attack.mitre.org/mitigations/M1022/>, 2020. Accessed: 2025-02-02.
- [210] MITRE ATT&CK, "Update Software." <https://attack.mitre.org/mitigations/M1051/>, Jul 2020. Accessed: 2025-02-02.
- [211] MITRE ATT&CK, "Pre-compromise." <https://attack.mitre.org/mitigations/M1056/>, Oct 2020. Accessed: 2025-02-02.
- [212] J. Mishra and S. K. Sahay, "Modern Hardware Security: A Review of Attacks and Countermeasures," 2025.
- [213] M. A. Butt, Z. Ajmal, Z. I. Khan, M. Idrees, and Y. Javed, "An in-depth survey of bypassing buffer overflow mitigation techniques," *Applied Sciences*, vol. 12, no. 13, p. 6702, 2022.
- [214] MITRE ATT&CK, "Behavior Prevention on Endpoint." <https://attack.mitre.org/mitigations/M1040/>, Jun 2019. Accessed: 2025-02-02.
- [215] MITRE ATT&CK, "Execution Prevention." <https://attack.mitre.org/mitigations/M1038/>, Oct 2024. Accessed: 2025-02-02.
- [216] MITRE ATT&CK, "Disable or Remove Feature or Program." <https://attack.mitre.org/mitigations/M1042/>, 2020. Accessed: 2025-02-02.
- [217] MITRE ATT&CK, "Limit Software Installation." <https://attack.mitre.org/mitigations/M1033/>, Oct 2024. Accessed: 2025-02-02.
- [218] MITRE ATLAS, "Use Multi-Modal Sensors." <https://atlas.mitre.org/mitigations/AML.M0009/>, Oct 2023. Accessed: 2025-02-01.
- [219] Q. Liu, M. Shoaib, M. U. Rehman, K. Bao, V. Hagenmeyer, and W. U. Hassan, "Accurate and Scalable Detection and Investigation of Cyber Persistence Threats," *arXiv preprint arXiv:2407.18832*, 2024.
- [220] A. Esmradi, D. W. Yip, and C. F. Chan, "A comprehensive survey of attack techniques, implementation, and mitigation strategies in large language models," in *International Conference on Ubiquitous Security*, pp. 76–95, Springer, 2023.
- [221] B. H. Zhang, B. Lemoine, and M. Mitchell, "Mitigating unwanted biases with adversarial learning," in *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, pp. 335–340, 2018.
- [222] R. R. Maiti, C. H. Yoong, V. R. Palleti, A. Silva, and C. M. Poskitt, "Mitigating adversarial attacks on data-driven invariant checkers for cyber-physical systems," *IEEE Transactions on Dependable and Secure Computing*, vol. 20, no. 4, pp. 3378–3391, 2022.
- [223] R. Huang and Y. Li, "Adversarial attack mitigation strategy for machine learning-based network attack detection model in power system," *IEEE Transactions on Smart Grid*, vol. 14, no. 3, pp. 2367–2376, 2022.
- [224] R. Koizumi and R. Sasaki, "Study on Countermeasures Using Mitigation Software against Vulnerability Attacks," in *2015 Fourth International Conference on Cyber Security, Cyber Warfare, and Digital Forensic (CyberSec)*, pp. 28–33, IEEE, 2015.
- [225] L. E. S. Jaramillo, "Malware detection and mitigation techniques: Lessons learned from Mirai DDOS attack," *Journal of Information Systems Engineering & Management*, vol. 3, no. 3, p. 19, 2018.
- [226] I. Branesco, O. Grigorescu, and M. Dascalu, "Automated Mapping of Common Vulnerabilities and Exposures to MITRE ATT&CK Tactics," *Information*, vol. 15, no. 4, p. 214, 2024.
- [227] N. Pinzón, V. Koundinya, R. E. Galt, W. O. Dowling, M. Baukloh, N. C. Taku-Forchu, T. Schohr, L. M. Roche, S. Ikendi, M. Cooper, et al., "AI-powered fraud and the erosion of online survey integrity: an analysis of 31 fraud detection strategies," *Frontiers in Research Metrics and Analytics*, vol. 9, p. 1432774, 2024.
- [228] MITRE ATLAS, "Sanitize Training Data." <https://atlas.mitre.org/mitigations/AML.M0007/>, Oct 2023. Accessed: 2025-02-01.
- [229] MITRE ATLAS, "Maintain AI Dataset Provenance." <https://atlas.mitre.org/mitigations/AML.M0025/>, Oct 2024. Accessed: 2025-02-01.
- [230] S. S. Balantrapu, "Adversarial Machine Learning: Security Threats and Mitigations," *International Journal of Sustainable Development in Computing Science*, vol. 1, no. 3, pp. 1–18, 2019.
- [231] MITRE ATLAS, "Validate ML Model." <https://atlas.mitre.org/mitigations/AML.M0008/>, Jan 2024. Accessed: 2025-02-01.
- [232] MITRE ATLAS, "Model Distribution Methods." <https://atlas.mitre.org/mitigations/AML.M0017/>, Jan 2024. Accessed: 2025-02-01.
- [233] MITRE ATLAS, "Control Access to ML Models and Data in Production." <https://atlas.mitre.org/mitigations/AML.M0019/>, Jan 2024. Accessed: 2025-02-01.
- [234] MITRE ATLAS, "AI Telemetry Logging." <https://atlas.mitre.org/mitigations/AML.M0024/>, Oct 2024. Accessed: 2025-02-01.
- [235] MITRE ATLAS, "Model Hardening." <https://atlas.mitre.org/mitigations/AML.M0003/>, Oct 2023. Accessed: 2025-02-01.
- [236] MITRE ATLAS, "Input Restoration." <https://atlas.mitre.org/mitigations/AML.M0010/>, Oct 2023. Accessed: 2025-02-01.
- [237] MITRE ATLAS, "Adversarial Input Detection." <https://atlas.mitre.org/mitigations/AML.M0015/>, Oct 2023. Accessed: 2025-02-01.
- [238] MITRE ATLAS, "Generative AI Guardrails." <https://atlas.mitre.org/mitigations/AML.M0020/>, Oct 2024. Accessed: 2025-02-01.
- [239] MITRE ATLAS, "Generative AI Guidelines." <https://atlas.mitre.org/mitigations/AML.M0021/>, Oct 2024. Accessed: 2025-02-01.
- [240] MITRE ATLAS, "Generative AI Model Alignment!" <https://atlas.mitre.org/mitigations/AML.M0022/>, Oct 2024. Accessed: 2025-02-01.
- [241] D. Lee and M. Tiwari, "Prompt Infection: LLM-to-LLM Prompt Injection within Multi-Agent Systems," 2024.
- [242] M. Ding, "Analyzing the Research and Application of the RSA Cryptosystem in the Context of Digital Signatures," *Transactions on Computer Science and Intelligent Systems Research*, vol. 5, Aug. 2024.
- [243] MITRE ATT&CK, "Application Isolation and Sandboxing." <https://attack.mitre.org/mitigations/M1048/>, 2020. Accessed: 2025-02-02.
- [244] MITRE ATT&CK, "Exploit Protection." <https://attack.mitre.org/mitigations/M1050/>, Jun 2020. Accessed: 2025-02-02.
- [245] MITRE ATT&CK, "Network Segmentation." <https://attack.mitre.org/mitigations/M1030/>, May 2020. Accessed: 2025-02-02.
- [246] MITRE ATT&CK, "Network Segmentation." <https://attack.mitre.org/mitigations/M1016/>, Jun 2020. Accessed: 2025-02-02.
- [247] E. Astranova and P. Issa, "AI-Powered Voice Spoofing for Next-Gen Vishing Attacks." <https://cloud.google.com/blog/topics/threat-intelligence/ai-powered-voice-spoofing-vishing-attacks>, Jul 2024. Accessed: 2025-01-29.
- [248] V. Cyber, "Can you trust ChatGPT's package recommendations?." <https://vulcan.io/blog/ai-hallucinations-package-risk>, Aug 2024. Accessed: 2025-01-26.
- [249] Google Threat Intelligence Group, "Adversarial Misuse of Generative AI." <https://cloud.google.com/blog/topics/threat-intelligence/adversarial-misuse-generative-ai>, Jan 2025. Accessed: 2025-01-29.
- [250] G. T. et al., "Gemma: Open Models Based on Gemini Research and Technology," 2024.
- [251] A. Vassilev, A. Oprea, A. Fordyce, and H. Anderson, "Adversarial Machine Learning: A Taxonomy and Terminology of Attacks and Mitigations," Tech. Rep. NIST AI 100-2e2023, National Institute of Standards and Technology, Gaithersburg, MD, 2024.
- [252] I. Kolochenko and G. Platt, "Data Security Professional Perspective: ChatGPT IP Cybersecurity." <https://www.bloomberglaw.com/external/document/X7K3GI38000000/data-security-professional-perspective-chatgpt-ip-cybersecurity->, Feb 2023. Accessed: 2025-01-29.