

A Survey on Privacy Attacks Against Digital Twin Systems in AI-Robotics

Ivan A. Fernandez*, Subash Neupane[†], Trisha Chakraborty[‡], Shaswata Mitra[§],
Sudip Mittal[¶], Nisha Pillai^{||}, Jingdao Chen^{**}, Shahram Rahimi^{††}

Computer Science and Engineering, Mississippi State University

Email: {*iaf28, [†]sn922, [‡]tc2006, [§]sm3843}@msstate.edu, {[¶]mittal, ^{||}pillai, ^{**}chenjingdao, ^{††}rahimi}@cse.msstate.edu

Abstract—Industry 4.0 has witnessed the rise of complex robots fueled by the integration of Artificial Intelligence/Machine Learning (AI/ML) and Digital Twin (DT) technologies. While these technologies offer numerous benefits, they also introduce potential privacy and security risks. This paper surveys privacy attacks targeting robots enabled by AI and DT models. Exfiltration and data leakage of ML models are discussed in addition to the potential extraction of models derived from first-principles (e.g., physics-based). We also discuss design considerations with DT-integrated robotics touching on the impact of ML model training, responsible AI and DT safeguards, data governance and ethical considerations on the effectiveness of these attacks. We advocate for a trusted autonomy approach, emphasizing the need to combine robotics, AI, and DT technologies with robust ethical frameworks and trustworthiness principles for secure and reliable AI robotic systems.

I. INTRODUCTION AND MOTIVATION

In today's interconnected world, the potential impact of cyberattacks on critical infrastructure cannot be overstated. Take, for example, a manufacturing facility brought to a standstill, not by a physical breakdown, but by a carefully orchestrated cyberattack. The target: the Digital Twin (DT) (virtual representations of physical systems) [1] of a critical robotic arm, which is an automated machine used for a complex assembly process. An adversary, exploiting vulnerabilities in the DT's Artificial Intelligence (AI) system, manipulates sensor readings, feeding false information about the robot's position and environment. The consequences of such a breach extend beyond operational disruptions (for example, production grinds to a halt, costly repairs, maintenance etc.); they undermine business operations and trust, posing serious risks to safety-critical environments where human lives might be at stake. This scenario, unfortunately, is not science fiction. With the increasing reliance on AI-enabled robotic DT systems across industries such as military [2], aerospace [3], intelligent manufacturing [4], healthcare [5], smart cities and transportation [1], [6], the security of their digital twins becomes paramount.

The convergence of Industry 4.0 technologies - AI, robotics, cloud computing, and the Internet-of-Things (IoT) - has given rise to sophisticated Cyber-Physical Systems (CPS). At the heart of this revolution are DTs. By leveraging AI and DT technologies, AI robots are now capable of higher levels of autonomy and complex Human-Robot Interactions (HRI) within multi-agent environments [7], finding applications in safety-

critical domains where human lives, system integrity, and environmental safety are at stake [8]. However, this reliance on AI and DTs introduces significant security risks, particularly in the realm of privacy. AI-enabled robotic DT models inherently rely on vast amounts of data, making them susceptible to breaches if inadequate data governance and cybersecurity measures are in place. An unsecured communication protocol, for instance, can provide unauthorized access to sensitive data and algorithms [9], [10], potentially compromising the entire system. The case of the Aethon TUG, a smart autonomous mobile robot used in hospitals, highlights this vulnerability. Researchers discovered critical vulnerabilities in the system that could allow adversaries to seize control and extract sensitive patient information [11].

This vulnerability is further amplified by the tight coupling between AI and DTs. Robotic DT models often incorporate AI models for analysis and decision-making, creating a complex ecosystem where a security flaw in one component can cascade into the other. Recognizing these risks, governments and regulatory bodies are developing guidelines for Responsible AI (RAI). The United States White House, for instance, issued an executive order in October 2023, establishing new standards for AI safety and security [12].

While the security risks of AI and DTs have been studied independently, the synergistic impact of their convergence on AI robotic systems, particularly from a privacy perspective, remains largely unexplored. This paper addresses this critical gap, investigating the unique challenges and vulnerabilities arising from the integration of AI and DTs in AI robots. To better understand and mitigate these emerging threats, we will utilize the MITRE Adversarial Threat Landscape for Artificial Intelligence Systems (ATLAS) [13] framework. While ATLAS provides a valuable resource for understanding AI-specific threats, we argue that it needs to be extended to adequately address the privacy risks associated with DTs. The main contributions of this paper are as follows:

- 1) Provide a survey of privacy attacks against DT systems used by AI robots.
- 2) Suggest the addition of first-principles models to the MITRE ATLAS framework under exfiltration.
- 3) Provide a discussion on how trusted autonomy can be achieved by combining robotics, AI, and DT technologies with ethics and trustworthiness concepts.

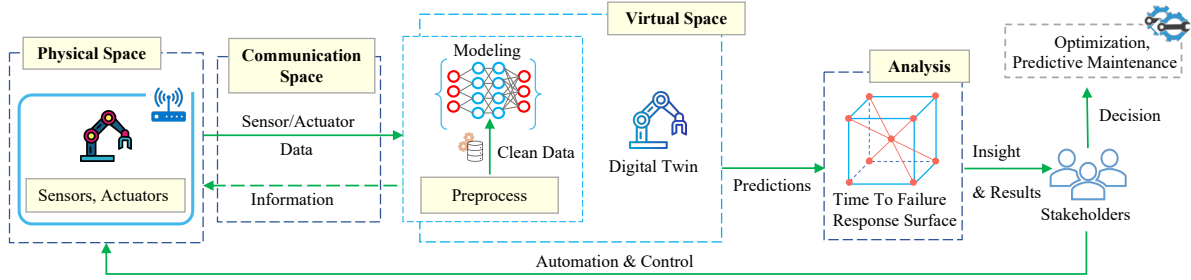


Fig. 1. Generic Framework of Data-Driven Robotic Digital Twin: Physical spaces comprise robotic sensors that collect data. Virtual space utilizes the data collected from physical space via a communication link between them. Predictions are generated by the AI models within virtual space, which are then analyzed before decisions are made by stakeholders.

The remainder of the paper investigates the security of AI and DT systems in robotics. In Section II we provide a primer on digital twin and then discuss its integration with robotics. Section III explores the attack surfaces on robotics systems that make of use AI and DT models. Following that, we touch on the impact of using machine learning (ML) models, responsible AI and DT safeguards, data governance and ethical consideration for DT-integrated robotics in Section IV. Finally, Section V concludes our paper.

II. BACKGROUND

In this section we provide the background on Digital Twin (DT) technology, and discuss the impact of digital twins to robotics including the importance effective data exchange.

A. Digital Twin Paradigm

The concept of using “physical twins”, which served as an early precursor to digital twins, has its historical origins in the 1970s, aligning with NASA’s Apollo mission [1]. In 2002, Grieves et al. [14] introduced the notion of DTs informally, later formalizing it in their published white paper. In 2012, NASA and United States Air Force (USAF) researchers [15] define a digital twin as “*an integrated multiphysics, multiscale, probabilistic simulation of an as-built vehicle or system that uses the best available physical models, sensor updates, fleet history, etc., to mirror the life of its corresponding flying twin.*” A simpler definition is that a DT is a virtual prototype of physical assets that simulates, emulates, mirrors, or twins the real-time operational conditions to behave like a real physical asset. A typical DT model as depicted in Fig. 1 comprises three components: a *physical space*, *virtual space*, and *communication space*. We explain the functions of each these components in greater details below.

1) *Physical Space*: The physical space is comprised of real-world objects, such as equipment, systems, cameras, sensors, components, etc. that are responsible for collecting data of the current physical measurement of objects. In the context of robotics systems, the two fundamentally connected physical components are sensors and actuators, which often work in tandem but are essentially opposite in their nature. For example, a sensor monitors the state and sends a signal when changes occur, whereas an actuator receives the signal and performs an action. Robots today are outfitted with several

sensors that generate volumes of operational data. These data are typically collected in time-series format, expressed as multi-channel sensor data, and stored in a datastore [16].

2) *Virtual Space*: In a DT system, virtual space can be viewed as a virtual replica that maintains a real-time model of its physical space. It receives the data collected in the physical space as input which is then preprocessed. The preprocessing step includes various sub-steps such as cleaning, removing redundant observations, transformation, data restructuring, and scaling. The output of the preprocessing step is a high-quality dataset. Deep learning models and/or first-principles models are then applied to this data in order to create the DT-based system. DT systems crafted in this manner have various capabilities. They can visualize the instant status of their physical space [17] through continuous monitoring, perform time-series forecasting on the remaining useful life [18] of a degradable vehicle component, detect abnormal patterns [1], and perform fault detection in data signals.

3) *Communication Space*: The communication space is the infrastructure that connects physical and virtual spaces [19]. It allows for information exchange between the various components of the overall DT ecosystem. Digital twins require effective communication to maintain synchronization and accurately interact with their real-world counterparts. Both wired (e.g., fiber optics, CAN Bus, ARINC-429) and wireless (e.g., WiFi, Zigbee, Bluetooth, 5G) may be used for data exchange.

4) *Analysis*: Once the digital twin is produced, the insights generated by it can be used to achieve the given robotic systems goals, such as repairs, prognostics, optimization, or predictive maintenance. In addition, the results can be utilized to update and improve both the physical and virtual twins, which will provide more autonomy and control.

B. Digital Twin-Integrated Robotics

An AI robot is an agent that can sense and act to maximize its chances of success for a given task [20]. Robots can have varying levels of autonomy (e.g., no autonomy, semi-autonomy, full-autonomy) based on how much human control is needed [21]. A fully autonomous AI robot is realizable by incorporating AI [22]. Furthermore, for an autonomous robot to be considered AI-enabled, it must demonstrate at least one of the fundamental constructs of AI such as *reasoning*, *planning*, *learning*, *communication*, and *perception* [23].

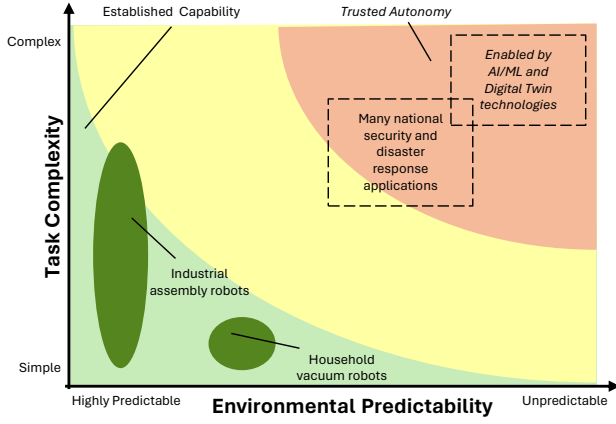


Fig. 2. Overview of autonomous capabilities adapted from Boulet et al. [24]. Two dimensions of the Autonomy Levels for Unmanned Systems (ALFUS) model [25] are shown. We make two key additions to the original: (1) “Future Capability” is replaced with “Trusted Autonomy” and (2) note that AI robots performing complex tasks in unpredictable environments must be enabled by some combination of AI/ML and DT technologies.

The concept of a DT is complementary methodology to AI [26], [27]. We argue that both AI and DT technologies are needed to achieve *Trusted Autonomy* as shown in Figure 2. AI aids with reducing the dimensionality of highly complex tasks and digital twin concepts facilitate exploration of different scenarios (or *What-If simulations* [28]) safely and efficiently.

DTs offer a powerful tool for designing, testing, operating, and maintaining robotic systems more effectively. They can be integrated with Model-Based Design (MBD) [29] and Model-Based Systems Engineering (MBSE) [30] processes to improve decision making throughout the robotic system’s lifecycle. Lie et al. [31] provide industrial applications for DTs in the design, development (manufacturing), service, and retirement phases of the lifecycle. During the design phase, DTs can be used for iterative optimization of the product design and support of virtual prototyping. Similarly, during the the manufacturing phase, the use of digital twins allows for real time monitoring and optimization of processes. Furthermore, DTs can be used for state monitoring, fault detection, predictive maintenance, and virtual testing during the service phase of the system.

DTs can be used to represent a robotic component, a robotic system, or a System of Systems (SoS). DTs can also be employed at the edge [32] for real-time decision making or in fog/cloud infrastructures if the use case allows for some latency. The rest of this section provides examples of DTs from different tiers with respect to the robotic system.

1) *Component-level DTs*: The DT paradigm allows for modeling of specific components in complex systems [33] like AI-enabled robots. Component-level DTs can be used as virtual models in lieu of physical hardware to test out requirements. Kutzke et al. [34] discuss how digital twins can be used for subsystems in Autonomous Underwater Vehicles (AUV) by establishing a generic process for determining a set of priority-based system components requiring digital

twin development for Condition-Based Maintenance (CBM) purposes. In order to create a testbed to optimize electric propulsion drive systems in autonomous vehicles, Rassolkin et al. [35] develop DTs using physical models and virtual sensors.

2) *System-level DTs*: Component-level DTs can be integrated and combined to model entire complex systems [36]. As an example, Xiong et al. [37] design and implement a simulation method for car-following scenarios of an autonomous vehicle using a digital twin for the secondary vehicle.

3) *SoS-level DTs*: Similar to developing a system DT with virtual models for its components, a System-of-Systems (SoS) DT can be constructed by aggregating multiple system-level models. DT modeling can be simple or highly complex, and Autonomous Mobile Robots (AMR) are usually modeled as complex SoS DTs [38].

Adept adversaries can target any and all of the DT tiers. It is important to note that while information may be isolated at a component-level, if that component is vital to operation of the system, then a successful exfiltration by an adversary compromises the entire the system. Similarly, if a system within a SoS becomes victim of an attack, then the entire SoS is potentially compromised as well.

III. ATTACKS ON DT-INTEGRATED AI ROBOTS

Digital twins provide an enabling technology for developing and safeguarding advanced robotic systems. As discussed in the previous section, digital twins allow for exploration of “What-If?” scenarios and analyses that may be prohibitive with the physical system since they can simulate the system’s behavior under different conditions (e.g., environmental, system configuration). The DT is comprised of software components (e.g., algorithms and models) [39] and unauthorized access to these digital assets could lead to devastating damage to stakeholders including the theft of Intellectual Property (IP), Private Personal Information (PPI), and reverse engineering.

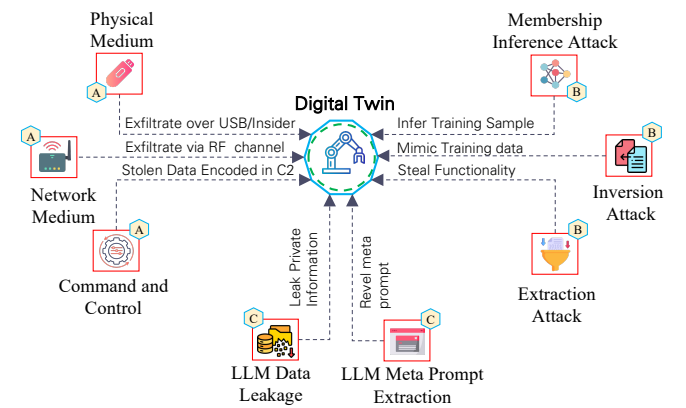


Fig. 3. A graphical illustration of different privacy exfiltration techniques. A represents exfiltration via cyber means, B represents exfiltration via Model inference API, while C represents possible attacks within LLM space.

The attack surface of a robotic system increases if DT or AI models are used. Adversaries could launch cyberattacks if the proper cybersecurity mechanisms (e.g., secured IoT networks

and protocols) are not in place for protecting these systems. The MITRE Adversarial Threat Landscape for Artificial-Intelligence Systems (ATLAS) [13] provides a knowledge base of tactics, techniques, and procedures (TTPs) against AI-enabled systems. There are 14 major tactics in the ATLAS but the main focus for this paper lies with *exfiltration*.

Malicious actors use privacy attacks to steal information about a physical system if they are able gain access to its DT. If the DT is comprised of a ML model, then an adversary could attempt to extract the model data (e.g., architecture, weights) or the training data (e.g., confidential data). This includes large-language models (LLMs) [40] if they are used by the AI robot. Privacy attacks can also occur against physics-based (first-principles) models. An adept adversary could query a physics-based model to generate enough input/target pairs to train a surrogate data-driven model. Privacy attacks are also known as confidentiality attacks because of their focus in the Confidentiality-Integrity-Availability (CIA) security triad.

Privacy attacks against DT-integrated AI robots could lead to catastrophic outcomes in safety-critical applications. This is because unauthorized access to a DT could lead to vulnerabilities in the physical twin (system) being leaked. Furthermore, allowing unauthorized access to sensitive information contained in a AI robot’s digital twin could lead to additional attacks (cyber or physical) and extortion.

Robot vacuums fall under the simple task complexity and highly predictable region in Figure 2. Their relatively known components and mass consumption can make them prime candidates for privacy attacks by adversaries. There are examples of these robots being used for eavesdropping via their LiDAR sensors [41]. In scenarios where AI robots must operate in unpredictable environments and perform a highly complex tasks, privacy attacks can lead to catastrophic damage. With the adoption of IoT concepts in the military domain [42] (i.e., Internet-of-Military-Things (IoMT)), autonomous military vehicles are also vulnerable to attacks. This was seen in December 2011, when Iranian hackers were able to bring down a RQ-170 Sentinel US spy drone [43].

The rest of this section will focus on privacy attacks against DT-integrated AI robots. In particular, we will discuss four types: *Exfiltration via Model Inference API*, *Exfiltration via Cyber Means*, *LLM Meta Prompt Extraction*, and *LLM Data Leakage* as depicted in Figure III. Table III-A provides a summary of these attacks.

A. Exfiltration via Model Inference API

Exfiltration is a form of data theft that occurs when an adversary steals artifacts and information about the AI robot from the DT ecosystem. Even though the use of Application Programming Interfaces (APIs) can to hide details of the model, exfiltration via model inference is still possible through black-box query access. In fact, this technique can be further divided into the following categories: *training data membership inference*, *model inversion*, and *model extraction*.

In Membership Inference Attacks (MIAs), adversaries aim to determine if a training sample belongs to the training

data of the targeted ML model. MIAs can target various ML models (e.g., classification, generative models [44]) and are usually constructed using one of two major approaches: binary classifier-based attack and metric-based attack [45]. In a binary classifier-based attack, *shadow training* [46] is usually used to train several shadow models as surrogate to the target model. Then an attack model (binary classifier) is trained to infer if a particular sample belonged to the target training data (member vs. non-member) from information provided by shadow models. Given a sample, metric-based inference attacks compares the output probability vector (confidence scores) provided by model against a predefined threshold as a way to determine if it belonged in the training data.

Expanding on metric-based inference attacks leads to label-only model inversion [47] where the goal is to infer sensitive information about the training data (or individual points) from output labels provided by the ML model. Han et al. [48] show how a model inversion attack against a black-box generative model can be constructed using its confidence scores as rewards to a Reinforcement Learning (RL) agent.

Model extraction is the process of creating a new model that approximates the behavior of the target model. Given enough queries, an adversary can careful craft training data to a surrogate ML model. The effect of this strategy can be amplified if the adversary how some knowledge (gray-box) of the foundational model if the target model has been fine-tuned [49]. Foundation models enable transfer learning [50] because they allow their weights to be adapted to new tasks via fine-tuning. The added benefit of using a pre-trained model also comes with risk because they can retain information about the original dataset [51].

Exfiltration is not restricted to just ML models and can also be used to approximate models tied to first-principles (e.g., physics-based). Similar to targeting ML models, a threat actor with API access to a digital twin that uses a first-principles model could launch a query-based attack to compile enough input-output training pairs to derive a data-driven surrogate model of the target. The use of ML to augment physics modeling is not a novel idea [41], [52]. To the best of our knowledge, this is the first research paper to discuss model extraction privacy attack of a physics-based model using a ML surrogate model.

B. Exfiltration via Cyber Means

A DT and its corresponding physical system can also be targeted by traditional cyberattacks. Once adversaries have performed reconnaissance to identify security gaps of a system, they can infiltrate the Robot-DT ecosystem and establish a footprint (e.g., command and control) either through network or physical access (e.g., insider threat). Adversaries may also attempt to exfiltrate data via a physical medium, as illustrated in Fig III, such as a removable drive [13] (e.g. external hard drive, USB drive, or other removable storage and processing device). Sensitive training data and models can exfiltrated through communication space channels, sometimes in chunks to avoid triggering network traffic thresholds and security

TABLE I
 PRIVACY ATTACKS AND EXFILTRATION TECHNIQUES AGAINST DT-INTEGRATED ROBOTIC SYSTEMS ADAPTED FROM MITRE ATLAS MODEL.
 (*) DENOTE PROPOSED ADDITIONS TO THE MITRE ATLAS MODEL.

Technique	DT Enabling Technology	Overview	Research
Exfiltration via Model Inference API	First-principles* (Physical, Mathematical, Statistical) or ML-enabled models	ML models are prone to leak sensitive information about its training data and adversaries can exfiltrate information through the model API. This can lead to private data leak or model extraction itself [53]. Non-ML models* are also vulnerable to exfiltration through surrogate creation using ML approaches.	[39], [45]–[49], [54], [55]
Exfiltration via Cyber Means	First-principles* (Physical, Mathematical, Statistical) or ML-enabled models	ML artifacts or other information, including non-ML models*, that may be relevant to adversaries’ goals could be exfiltrated through traditional cyber means [56].	[39], [54], [57]–[59]
LLM Meta Prompt Extraction	LLM-enabled models	An adversary may induce an LLM to disclose its internal instruction prompts, also called “meta prompt.” The leak of meta prompt can equip an adversary with the knowledge of the system’s internal workings, and security policies or even may lead to intellectual property theft [60].	[61]–[66]
LLM Data Leakage	LLM-enabled models	An adversary may craft prompts purposely to induce an LLM to disclose sensitive information. Sensitive information can include private user data or proprietary information and may come from proprietary training data, retrieval sources the LLM is connected to, or information from other users of the LLM [67].	[40], [61], [68]–[73]

mechanisms. This is a well studied research area and readers are directed to the MITRE ATT&CK framework [74] and recent papers for more insight [75]–[77].

An important addition that this paper provides to the literature is that models based on first principles (i.e., non-ML models) are also susceptible to this attack strategy. Borky et al. [78] advice on the use of cybersecurity practices within a MBSE process in order to protect sensitive information. MBSE and physics-based models tie-back to requirements [79] and known behavior [80] for the system. Depending on the motivation of the adversary, successful exfiltration of sensitive information allows for extortion for an immediate gain. It also makes the system vulnerable to future attacks for the rest of its lifecycle if the initial attack goes unnoticed and the discovered vulnerabilities go unmitigated.

C. LLM Meta Prompt Extraction

DT-integrated AI robots can leverage Large Language Models (LLMs) to augment the Human-Robot Interaction (HRI) [81]. LLMs are well known for their few-shot learning [82] and In-Context Learning (ICL) capabilities which makes them widely popular and readily accessible to non AI/ML experts. Adept adversaries can attempt to circumvent the safety constraints and safeguarding mechanisms [64] of a targeted LLM to retrieve its initial instructions (i.e., meta prompt). This technique is sometimes called a “jailbreaking” attack [65] because it allows the adversary bypass guardrails to access internal workings of the model through strategic inputs such as role-playing prompts [66]. With the knowledge gained through LLM meta prompt extraction, an adversary can launch additional attacks including theft of sensitive information. If the LLM is a crucial component to the operation and decision-making process of an AI robot, then a successful jailbreaking attack could lead to a compromised system.

D. LLM Data Leakage

As described earlier, LLMs are vulnerable to leaking sensitive data [40] if special care is not taken by integrators. Studies have shown that LLMs are capable of leaking their training data [71] due to memorization [72]. Besides training

data, LLMs are also capable of leaking connected data sources, information from other users, and model information. Similar to black-box extraction attacks to other ML models,

Birch et al. [73] show a “model leeching” attack that is cost-effective and capable of distilling task-specific knowledge from a target LLM to generate a reduced parameter model. Model leeching is described as a black-box adversarial attack that aims to extract the target LLM by creating a copy (i.e., recreating the target model) within a specific task. To recreate the model, three stages are needed: prompt design, data generation, and (stolen) model training. If successful model recreation occurs, then adversaries can use it as a surrogate in a staging ground to launch additional attacks against the target model.

Finlayson et al. [83] show that it is possible to learn a significant amount of non-public information from API-protected LLMs with a small amount of queries. The approach exploits the low-rank output layer common to most LLM architectures. In their research, it shown that the softmax bottleneck imposes low-rank constraints on LLM outputs. They leverage the restricted output space to obtain the LLM image using a small number of LLM outputs. The LLM image can be thought of as a model signature and it exposes model hyperparameters, output layer parameters, and full model outputs.

IV. DT-INTEGRATED ROBOTICS DESIGN CONSIDERATIONS AND DISCUSSION

In previous sections, we discussed how AI and DT technologies are essential for enabling robots to perform complex tasks in unpredictable environments including safety-critical applications where human-life could be at risk. These technologies can be used in any and all stages of the the AI robot’s lifecycle from design (beginning of life) to retirement (end of life) [84]. The use of these technologies also provides a pathway for threat actors to carry out cyberattacks including those that aim to extract private information about the system.

In this section, we touch on the impact of ML model training, responsible AI/DT safeguards, and ethical considerations to the effectiveness of privacy attacks.

A. Risk Factors in ML-enabled Digital Twins

If a ML-enabled digital twin is used during any phase of an AI robot's lifecycle, special care must be taken in securing those models because factors such as overfitting, dataset structure, model architecture, and model type could affect the accuracy of certain privacy attacks [85], [86].

1) *ML Model Overfitting and Data Structure*: Overfitting occurs when the underlying ML model in a data-driven DT believes the noise and outliers of the training data to the extent that it performs poorly on new, unseen data. By memorizing details of the training instances, these ML models may not learn the underlying trends in the data and could leave them at risk for privacy attacks. In their research, Yeom et al. [87] introduce a series of formal definitions for examining the effects of overfitting to membership inference and attribute inference attacks. The theoretical and experimental results of their study show that ML models are more vulnerable to attacks as more overfitting occurs. As part of their study, they analyze linear regression, tree, and Convolutional Neural Network (CNN) models on different datasets (e.g., IWPC, Netflix). Yeom et al. [88] later show that overfitting is not a necessary condition to privacy attacks and that models trained to be robust against adversarial examples are also exposed. They note that defending against both privacy and integrity attacks simultaneously may be challenging in some cases.

The composition of the training data can contribute the vulnerability of the ML-enabled DT to privacy attacks. Parallel to overfitting is the notion that if a model has to memorize out-of-distribution samples (i.e., outliers) found in the training data, then those samples are vulnerable to attacks. In addition, the use Personal Identifiers (PIDs) as part of the model training or inference only increases the privacy attack surface [89].

2) *ML Model Complexity*: Similar to overfitting, model complexity can also affect data privacy. Highly complex models, especially deep learning models, contain large number of parameters that could lead them to memorizing the training data instead of learning patterns for generalization. Another added effect is that complex models are usually obscure in terms of interpretability which make finding potential data leaks difficult.

In order to address transparency and to provide the user a level of confidence, AI practitioners may augment their models with explainable AI (XAI) techniques. Special care must be taken to balance both explainability and privacy because the use XAI techniques can lead to privacy risks [90]. With access to explanations, adversaries could craft stronger privacy attacks including member inference and model inversion. Shokri et al. [91] show that backpropagation-based techniques are capable of leaking membership information because of their ability to statistically characterize decision boundaries. Zhao et al. [92] develop an image-based, XAI-aided inversion model with emotion prediction as the target task and face reconstruction as the attack task.

The bulk of the privacy attack research is focused on generative models [93], [94] including Variational Autoencoders and Generative Adversarial Networks. Previously, we

discussed how LLMs are also vulnerable to leaking sensitive information. Today, the largest and most capable LLMs are using a transformer-based architecture [95]. Due to their popularity, there's a rising trend to use transformers for other tasks. Lu et al. [96] introduce the Attention Privacy Leakage (APRIL) attack to steal private local training data from shared gradients of a Vision Transformer (ViT). The attack showcases the vulnerability of learnable position embeddings.

B. Preventing Privacy Attacks

While the aim of this paper is to provide an overview of privacy attacks against DT-integrated AI robots, it is also important to discuss defensive strategies against these attacks. There are several defenses against privacy attacks (or Privacy Enhancing Techniques (PET)) including anonymization, encryption, and differential privacy [97]. The underlying principle behind these defenses are tied to *The Fundamental Law of Information Recovery* [98] which formulates the notion that "giving overly accurate answers to too many questions will inevitably destroy privacy." The law was formalized and proven with *reconstruction attacks* [99] which describes any method that aggregates (i.e., compiles) publicly-available information to partially recreate a private dataset. Next, we will describe several major defensive strategies for the privacy attack techniques described in Table III-A.

For query-based privacy attacks, one simple mitigating approach is to restrict the number of model queries. This prevents adversaries from compiling enough question-answer (QA) pairs to carry out their attack. *Anonymization* techniques can also be employed to remove personable identifiers and information [97] from the training data. This strategy is usually not useful on its own since there are de-anonymization approaches available to adversaries [100].

Homomorphic encryption [101] is another defensive strategy against this type of attacks since it allows for computations to be performed on encrypted data. With Homomorphic Encryption (HE), only the user with the matching private key is able to decrypt data to reveal its contents. A practical HE approach has been shown with transformers but it comes with a performance burden [102].

Differential privacy is closely tied to reconstruction attacks and *The Fundamental Law of Information Recovery* [98]. Differential privacy (DP) can be achieved by adding noise to the training data, model parameters, or model outputs [103]. Recently, DP has been demonstrated on LLMs [104].

Model watermarking [105] could also be used in order to protect a model's Intellectual Property (IP) rights. This approach embeds unique identifiers into ML models allowing for traceability back to the model if it is ever falls victim to an exfiltration attack. The watermarking strategy would have to be robust to "watermark overwriting" attacks that adversaries could employ [106].

In order to secure DTs used in robotics, it is crucial that cybersecurity best practices are used at the onset of their lifecycle. This means that the information and artifacts used

from design to deployment are protected using strict access controls and authentication mechanisms.

C. Trustworthy, Ethical, and Responsible DT

Data-driven DTs generate predictions, which are then analyzed to extract insights. Stakeholders subsequently use these insights to make decisions, such as whether to perform maintenance or optimize a component. However, the AI models utilized for prediction are typically black-box models, lacking clear explainability. In other words, they are unable to justify their judgments and predictions [107]. A trustworthy DT should enable users to comprehend the decision-making process and the rationale behind its actions. One way to enhance trustworthiness in DTs that leverage AI models is to ensure the explainability of these models. This can be achieved through two approaches: employing intrinsically transparent models for prediction, such as decision trees or linear models, or generating explanations [1] after a decision has been made, also known as post-hoc explanations [108].

Apart from predictions, AI robots must be trusted to behave safely and ethically. For instance, DTs of autonomous vehicle's (AV) perception systems must have sufficient situational awareness, for it to make the right decisions, to keep itself and humans safe [109]. Isaac Asimov's three laws of robotics [110] can be used as guidelines for programming AI agents within robotic systems to augment trust.

Ethics has various definitions and approaches in the literature. Kuipers [109] defines ethics as a "*set of beliefs that a society conveys to its individual members, to encourage them to engage in positive-sum interactions and to avoid negative-sum interactions.*" Robotic systems and applications are governed by *roboethics* [111], which oversees the ethical outcomes and aftermaths arising from robotics technology. When it comes to artificial intelligence, AI ethics outlines the moral responsibilities and duties of both the AI system and its developers. The combination of roboethics and AI ethics plays a pivotal role in shaping the development and deployment of responsible DT-Integrated AI robots. Responsible DT can be achieved by considering three levels of AI ethics when designing intelligent systems: *ethics in design*, *ethics by design*, and *ethics for design* [112]. The first approach ensures that the development process takes into account the ethical and societal implications of AI and its role in the socio-technical environment. The second approach is concerned with integrating ethical reasoning abilities as part of the behavior of artificial autonomous systems. The third approach, on the other hand, focuses on the research integrity of stakeholders (researchers, developers) and institutions to ensure regulation and certification. Overall, a responsible DT is ethical, lawful, reliable, and beneficial.

D. Data Governance and Data Management

AI-enabled digital twins, including those used in robotics, must be managed and protected using *data governance* practices and processes. In fact, data governance (DG) lays the foundation for Trustworthy AI [113]. The purpose of DG is to

build trust in data [114]. It encompasses the policies, standards, and practices required to effectively manage data throughout its lifecycle. The use of data governance frameworks ensures that data is accurate, available, and secure.

Transparency is crucial in order to build trust and accountability with stakeholders. Important questions that need to be addressed in data-driven processes include: *What are the types of data being collected? How will the data be used? Who will have access to the data?* Organizations involved in DT typically establish contractual agreements and develop data governance frameworks to clarify *data ownership*. Compliance with data protection regulations, such as GDPR [115] or the CCPA [116], may require data ownership declarations and obtaining consent for data collection/usage in order to ensure transparency.

Tied to transparency is the concept of traceability [117]. The data collected by DTs, processed by (e.g., features) DTs, and made by DTs (e.g., decisions, inferences) must be logged to allow for traceability. The data transformations that occur in AI robots throughout data-driven processes are operating on behalf of individuals or making autonomous decisions. For these decisions to be ethically and responsibly taken, they must be explainable and align with transparency requirements which ensures fairness and democracy [118]–[120]. This is especially true in safety-critical applications where loss of life or the environment may be damaged if the AI robot fails. *Who is responsible when AI robots fail? How can we accurately determine what went wrong when the system fails?*

E. Trusted Autonomy

Trusted autonomy refers to the design, development, and deployment of AI robots that can be relied upon to operate safely, effectively, and ethically within their intended environments. Experts [121] describe research in trusted autonomy as covering multiple "advanced topics such as robotics, AI, simulation and ethics." Despite significant progress, achieving trusted autonomy is an ongoing effort that requires continuous improvement and adaptation to new challenges.

We believe that trusted autonomy can be achieved by combining robotics with AI and DT technologies. In addition, ethics and trustworthiness procedures and cybersecurity processes are needed. AI allows for efficiently searching high-dimensional spaces [122] and pattern recognition. The latter allows for Dimensionality Reduction (DR) [123] and data compression [124]. It is hard for humans to visualize data beyond a low number of dimensions (e.g., 3D). Existing techniques along with the use topological AI and Topological Data Analysis (TDA) [125] can help reduce a high dimensional problem space. In doing so, it improves the explainability of these data-driven approaches allowing stakeholders to better visualize why decisions are being made.

AI and digital twin models can be integrated with Modeling and Simulation (M&S) frameworks allowing for "What-If" scenarios. Stakeholders can use make use of a digital twin in lieu of the physical system to explore and analyze how the asset would act in simulated environmental configurations.

In other words, the use of a DT allows for Testing and Evaluation (T&E) of the physical system if the models provide the needed fidelity. This is extremely important in safety-critical applications where human-life, the environment, or the system must be protected.

Data governance is the connective tissue that ties accuracy, availability, and security of data in the data-driven models and processes used by AI robots. Cybersecurity principles must also be used to ensure that models are secure and protected from adversaries looking to do harm (e.g., obtain unauthorized access to sensitive information). The merger of all the previous discussed concepts leads us into responsible and ethical use of AI robots and closer to trusted autonomy.

V. CONCLUSION

The convergence of AI and DT technologies presents both opportunities and challenges for the development of complex robots. As these robots become increasingly integrated into safety-critical applications, ensuring their security and trustworthiness is paramount. As illustrated in our survey, both ML-enabled and physics-based robotic DT models are susceptible to various privacy attacks aimed at extracting sensitive information, potentially compromising not only intellectual property but also human safety. To address these challenges, we advocate for a holistic approach that integrates cybersecurity best practices, robust privacy-preserving techniques, and adherence to ethical principles throughout the robot's DT lifecycle. Furthermore, we highlight the importance of the trusted autonomy, where the design and operation of these robots prioritize transparency, explainability, and rigorous security measures, ensuring they operate safely, responsibly, and transparently. The development and deployment of trustworthy, ethical, and secure DT-integrated AI robots will be crucial for their successful adoption in various sectors.

ACKNOWLEDGEMENTS

This work was supported by the Predictive Analytics and Technology Integration (PATENT) Laboratory at the Department of Computer Science and Engineering, Mississippi State University. The authors would like to acknowledge Dr. Jorge A. O'Farrill and Dr. Stephen R. Snarski (Technical Fellows, Modern Technology Solutions, Inc.) for their contributions to both Digital Twin and Trusted Autonomy topics.

REFERENCES

- [1] S. Neupane, I. A. Fernandez, W. Patterson, S. Mittal, M. Parmar, and S. Rahimi, "Twinexplainer: Explaining predictions of an automotive digital twin," *preprint arXiv:2302.00152*, 2023.
- [2] J. G. Metcalf, J. A. Laffey, and G. R. Cook, "Integrating digital twin concepts to enhance agility of the united states marine corps' decision support framework," 2023.
- [3] M. Xiong, H. Wang, Q. Fu, and Y. Xu, "Digital twin-driven aero-engine intelligent predictive maintenance," *The International Journal of Advanced Manufacturing Technology*, vol. 114, no. 11, 2021.
- [4] Y. Xu, Y. Sun, X. Liu, and Y. Zheng, "A digital-twin-assisted fault diagnosis using deep transfer learning," *Ieee Access*, vol. 7, 2019.
- [5] A. Croatti, M. Gabellini, S. Montagna, and A. Ricci, "On the integration of agents and digital twins in healthcare," *Journal of Medical Systems*, vol. 44, no. 9, p. 161, 2020.
- [6] O. El Marai, T. Taleb, and J. Song, "Roads infrastructure digital twin: A step toward smarter cities realization," *IEEE Network*, vol. 35, no. 2, pp. 136–143, 2020.
- [7] S. Neupane, S. Mitra, I. A. Fernandez, S. Saha, S. Mittal, J. Chen, N. Pillai, and S. Rahimi, "Security considerations in ai-robotics: A survey of current methods, challenges, and opportunities," *IEEE Access*, 2024.
- [8] A. Mazumder, M. Sahed, Z. Tasneem, P. Das, F. Badal, M. Ali, M. Ahamed, S. Abhi, S. Sarker, S. Das *et al.*, "Towards next generation digital twin in robotics: Trends, scopes, challenges, and future," *Heliyon*, 2023.
- [9] A. Botta, S. Rotbei, S. Zinno, and G. Ventre, "Cyber security of robots: A comprehensive survey," *Intelligent Systems with Applications*, 2023.
- [10] C. Cerrudo and L. Apa, "Hacking robots before skynet," *IOActive Website*, pp. 1–17, 2017.
- [11] Cynerio, "Jekyllbot 5 vulnerability disclosure report," <https://www.cynerio.com/jekyllbot-5-vulnerability-disclosure-report>, 2022, accessed: 2024-05-01.
- [12] J. R. Biden, "Executive order on the safe, secure, and trustworthy development and use of artificial intelligence," 2023.
- [13] M. ATLAS, "Adversarial threat landscape for artificial-intelligence systems," <https://atlas.mitre.org>, 2024, accessed: 2024-03-01.
- [14] M. Grieves, "Digital twin: manufacturing excellence through virtual factory replication," *White paper*, vol. 1, no. 2014, pp. 1–7, 2014.
- [15] E. Glaessgen and D. Stargel, "The digital twin paradigm for future nasa and us air force vehicles," in *53rd AIAA/ASME/ASCE/AHS/ASC structures, structural dynamics and materials conference 20th AIAA/ASME/AHS adaptive structures conference 14th AIAA*, 2012.
- [16] S. Neupane, I. A. Fernandez, W. Patterson, S. Mittal, and S. Rahimi, "A temporal anomaly detection system for vehicles utilizing functional working groups and sensor channels," in *2022 IEEE 8th International Conference on Collaboration and Internet Computing (CIC)*. IEEE, 2022, pp. 99–108.
- [17] T. H. Luan, R. Liu, L. Gao, R. Li, and H. Zhou, "The paradigm of digital twin communications," *arXiv preprint arXiv:2105.07182*, 2021.
- [18] K. Feng, J. Ji, Y. Zhang, Q. Ni, Z. Liu, and M. Beer, "Digital twin-driven intelligent assessment of gear surface degradation," *Mechanical Systems and Signal Processing*, vol. 186, p. 109896, 2023.
- [19] D. M. Botín-Sanabria, A.-S. Mihaita, R. E. Peimbert-García, M. A. Ramírez-Moreno, R. A. Ramírez-Mendoza, and J. d. J. Lozoya-Santos, "Digital twin technology challenges and applications: A comprehensive review," *Remote Sensing*, vol. 14, no. 6, p. 1335, 2022.
- [20] R. R. Murphy, *Introduction to AI robotics*. MIT press, 2019.
- [21] N. Melenbrink, J. Werfel, and A. Menges, "On-site autonomous construction robots: Towards unsupervised building," *Automation in construction*, vol. 119, p. 103312, 2020.
- [22] S. Ness, N. J. Shepherd, and T. R. Xuan, "Synergy between ai and robotics: A comprehensive integration," *Asian Journal of Research in Computer Science*, vol. 16, no. 4, pp. 80–94, 2023.
- [23] S. Samoil, M. L. Cobo, E. Gómez, G. De Prato, F. Martínez-Plumed, and B. Delipetrev, "Ai watch. defining artificial intelligence. towards an operational definition and taxonomy of artificial intelligence," 2020.
- [24] M. T. Boulet, "The autonomous systems tidal wave," *Lincoln Laboratory Journal*, vol. 22, no. 2, p. 19, 2017.
- [25] H.-M. Huang, K. Pavak, B. Novak, J. Albus, and E. Messin, "A framework for autonomy levels for unmanned systems (alfus)," *Proceedings of the AUVSI's unmanned systems North America*, pp. 849–863, 2005.
- [26] A. Fuller, Z. Fan, C. Day, and C. Barlow, "Digital twin: Enabling technologies, challenges and open research," *IEEE access*.
- [27] M. M. Rathore, S. A. Shah, D. Shukla, E. Bentafat, and S. Bakiras, "The role of ai, machine learning, and big data in digital twinning: A systematic literature review, challenges, and opportunities," *IEEE Access*, vol. 9, pp. 32 030–32 052, 2021.
- [28] F. Pires, B. Ahmad, A. P. Moreira, and P. Leitão, "Digital twin based what-if simulation for energy management," in *2021 4th IEEE International Conference on Industrial Cyber-Physical Systems (ICPS)*. IEEE, 2021, pp. 309–314.
- [29] G. Bachelor, E. Brusa, D. Ferretto, and A. Mitschke, "Model-based design of complex aeronautical systems through digital twin and thread concepts," *IEEE Systems Journal*, 2019.
- [30] A. M. Madni, C. C. Madni, and S. D. Lucero, "Leveraging digital twin technology in model-based systems engineering," *Systems*, 2019.

- [31] M. Liu, S. Fang, H. Dong, and C. Xu, "Review of digital twin about concepts, technologies, and industrial applications," *Journal of Manufacturing Systems*, vol. 58, pp. 346–361, 2021.
- [32] L. Girletti, M. Groshev, C. Guimarães, C. J. Bernardos, and A. de la Oliva, "An intelligent edge-based digital twin for robotics," in *2020 IEEE Globecom Workshops (GC Wkshps)*.
- [33] X. Zheng, J. Lu, and D. Kiritis, "The emergence of cognitive digital twin: vision, challenges and opportunities," *International Journal of Production Research*, vol. 60, no. 24, pp. 7610–7632, 2022.
- [34] D. T. Kutzke, J. B. Carter, and B. T. Hartman, "Subsystem selection for digital twin development: A case study on an unmanned underwater vehicle," *Ocean Engineering*, 2021.
- [35] A. Rassölkin, T. Vaimann, A. Kallaste, and V. Kuts, "Digital twin for propulsion drive of autonomous electric vehicle," in *2019 IEEE 60th International Scientific Conference on Power and Electrical Engineering of Riga Technical University (RTUCon)*, 2019.
- [36] W. Jia, W. Wang, and Z. Zhang, "From simple digital twin to complex digital twin part ii: multi-scenario applications of digital twin shop floor," *Advanced Engineering Informatics*, vol. 56, p. 101915, 2023.
- [37] H. Xiong, Z. Wang, G. Wu, Y. Pan *et al.*, "Design and implementation of digital twin-assisted simulation method for autonomous vehicle in car-following scenario," *Sensors*, 2022.
- [38] P. Stączek, J. Pizoń, W. Danilczuk, and A. Gola, "A digital twin approach for the improvement of an autonomous mobile robots (amr's) operating environment—a case study," *Sensors*, 2021.
- [39] C. Alcaraz and J. Lopez, "Digital twin: A comprehensive survey of security threats," *IEEE Communications Surveys & Tutorials*, vol. 24, no. 3, pp. 1475–1503, 2022.
- [40] J. Evertz, M. Chlosta, L. Schönherr, and T. Eisenhofer, "Whispers in the machine: Confidentiality in llm-integrated systems," *arXiv preprint arXiv:2402.06922*, 2024.
- [41] T. Zubatiuk and O. Isayev, "Development of multimodal machine learning potentials: toward a physics-aware artificial intelligence," *Accounts of Chemical Research*, vol. 54, no. 7, pp. 1575–1585, 2021.
- [42] M. Pradhan and J. Noll, "Security, privacy, and dependability evaluation in verification and validation life cycles for military iot systems," *IEEE Communications Magazine*, vol. 58, no. 8, pp. 14–20, 2020.
- [43] M. Yahuza, M. Y. I. Idris, I. B. Ahmedy, A. W. A. Wahab, T. Nandy, N. M. Noor, and A. Bala, "Internet of drones security and privacy issues: Taxonomy and open challenges," *IEEE Access*, vol. 9, pp. 57 243–57 270, 2021.
- [44] Y. Zhang, R. Jia, H. Pei, W. Wang, B. Li, and D. Song, "The secret revealer: Generative model-inversion attacks against deep neural networks," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 253–261.
- [45] H. Hu, Z. Salicic, L. Sun, G. Dobbie, P. S. Yu, and X. Zhang, "Membership inference attacks on machine learning: A survey," *ACM Computing Surveys (CSUR)*, vol. 54, no. 11s, pp. 1–37, 2022.
- [46] R. Shokri, M. Stronati, C. Song, and V. Shmatikov, "Membership inference attacks against machine learning models," in *2017 IEEE symposium on security and privacy (SP)*. IEEE, 2017.
- [47] T. Zhu, D. Ye, S. Zhou, B. Liu, and W. Zhou, "Label-only model inversion attacks: Attack with the least information," *IEEE Transactions on Information Forensics and Security*, 2022.
- [48] G. Han, J. Choi, H. Lee, and J. Kim, "Reinforcement learning-based black-box model inversion attacks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023.
- [49] Z. Li, J. Hong, B. Li, and Z. Wang, "Shake to leak: Fine-tuning diffusion models can amplify the generative privacy risk," *arXiv preprint arXiv:2403.09450*, 2024.
- [50] R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, M. S. Bernstein, J. Bohg, A. Bosselut, E. Brunskill *et al.*, "On the opportunities and risks of foundation models," *arXiv preprint arXiv:2108.07258*, 2021.
- [51] J. H. Lee and S. J. You, "Balancing privacy and accuracy: Exploring the impact of data anonymization on deep learning models in computer vision," *IEEE Access*, 2024.
- [52] J. Willard, X. Jia, S. Xu, M. Steinbach, and V. Kumar, "Integrating physics-based modeling with machine learning: A survey," *arXiv preprint arXiv:2003.04919*, vol. 1, no. 1, pp. 1–34, 2020.
- [53] MITRE, "Exfiltration via ml inference api," <https://atlas.mitre.org/techniques/AML.T0024>, 2024, accessed: 2024-03-01.
- [54] D. Holmes, M. Papathanasaki, L. Maglaras, M. A. Ferrag, S. Nepal, and H. Janicke, "Digital twins and cyber security—solution or challenge?" in *2021 6th South-East Europe Design Automation, Computer Engineering, Computer Networks and Social Media Conference (SEEDA-CECNSM)*, 2021.
- [55] J. Scheibmeir and Y. Malaiya, "An api development model for digital twins," in *2019 IEEE 19th International Conference on Software Quality, Reliability and Security Companion (QRS-C)*. IEEE, 2019.
- [56] MITRE, "Exfiltration via cyber means," <https://atlas.mitre.org/techniques/AML.T0025>, 2024, accessed: 2024-03-01.
- [57] J. E. Rubio, R. Roman, C. Alcaraz, and Y. Zhang, "Tracking apts in industrial ecosystems: A proof of concept," *Journal of Computer Security*, vol. 27, no. 5, pp. 521–546, 2019.
- [58] A. Saad, S. Faddel, T. Youssef, and O. A. Mohammed, "On the implementation of iot-based digital twin for networked microgrids resiliency against cyber attacks," *IEEE transactions on smart grid*, vol. 11, no. 6, pp. 5138–5150, 2020.
- [59] M. Eckhart and A. Ekelhart, "Towards security-aware virtual environments for digital twins," in *Proceedings of the 4th ACM workshop on cyber-physical system security*, 2018, pp. 61–72.
- [60] MITRE, "Llm meta prompt extraction," <https://atlas.mitre.org/techniques/AML.T0056>, 2024, accessed: 2024-03-01.
- [61] K. Greshake, S. Abdelnabi, S. Mishra, C. Endres, T. Holz, and M. Fritz, "Not what you've signed up for: Compromising real-world llm-integrated applications with indirect prompt injection," in *Proceedings of the 16th ACM Workshop on Artificial Intelligence and Security*, 2023, pp. 79–90.
- [62] Z. Sha and Y. Zhang, "Prompt stealing attacks against large language models," *arXiv preprint arXiv:2402.12959*, 2024.
- [63] Y. Yang, X. Zhang, Y. Jiang, X. Chen, H. Wang, S. Ji, and Z. Wang, "Prsa: Prompt reverse stealing attacks against large language models," *arXiv preprint arXiv:2402.19200*, 2024.
- [64] F. Wu, N. Zhang, S. Jha, P. McDaniel, and C. Xiao, "A new era in llm security: Exploring security concerns in real-world llm-based systems," *preprint arXiv:2402.18649*, 2024.
- [65] A. Wei, N. Haghtalab, and J. Steinhardt, "Jailbroken: How does llm safety training fail?" *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [66] J. Chu, Y. Liu, Z. Yang, X. Shen, M. Backes, and Y. Zhang, "Comprehensive assessment of jailbreak attacks against llms," *arXiv preprint arXiv:2402.05668*, 2024.
- [67] MITRE, "Llm data leakage," atlas.mitre.org/techniques/AML.T0057, 2024, accessed: 2024-03-01.
- [68] Y. Liu, G. Deng, Y. Li, K. Wang, T. Zhang, Y. Liu, H. Wang, Y. Zheng, and Y. Liu, "Prompt injection attack against llm-integrated applications," *arXiv preprint arXiv:2306.05499*, 2023.
- [69] A. Namer, J. Miller, H. Vagts, and B. Maltzman, "A cost-effective method to prevent data exfiltration from llm prompt responses," 2023.
- [70] S. S. Kumar, M. Cummings, and A. Stimpson, "Strengthening llm trust boundaries: A survey of prompt injection attacks."
- [71] N. Carlini, F. Tramer, E. Wallace, M. Jagielski, A. Herbert-Voss, K. Lee, A. Roberts, T. Brown, D. Song, U. Erlingsson *et al.*, "Extracting training data from large language models," in *30th USENIX Security Symposium (USENIX Security 21)*, 2021, pp. 2633–2650.
- [72] J. Huang, H. Shao, and K. C.-C. Chang, "Are large pre-trained language models leaking your personal information?" *arXiv preprint arXiv:2205.12628*, 2022.
- [73] L. Birch, W. Hackett, S. Trawicki, N. Suri, and P. Garraghan, "Model leeching: An extraction attack targeting llms," *arXiv preprint arXiv:2309.10544*, 2023.
- [74] B. E. Strom, A. Applebaum, D. P. Miller, K. C. Nickels, A. G. Pennington, and C. B. Thomas, "Mitre att&ck: Design and philosophy," in *Technical report*. The MITRE Corporation, 2018.
- [75] R. Lanotte, M. Merro, A. Munteanu, and L. Viganò, "A formal approach to physics-based attacks in cyber-physical systems," *ACM Transactions on Privacy and Security (TOPS)*, vol. 23, no. 1, 2020.
- [76] R. Alguliyev, Y. Imamverdiyev, and L. Sukhostat, "Cyber-physical systems and their security issues," *Computers in Industry*, vol. 100, pp. 212–223, 2018.
- [77] S. Neupane, I. A. Fernandez, S. Mittal, and S. Rahimi, "Impacts and risk of generative ai technology on cyber defense," *arXiv preprint arXiv:2306.13033*, 2023.

- [78] J. M. Borky, T. H. Bradley, J. M. Borky, and T. H. Bradley, "Protecting information with cybersecurity," *Effective Model-Based Systems Engineering*, pp. 345–404, 2019.
- [79] A.-L. Bruggeman, B. van Manen, T. van der Laan, T. van den Berg, and G. La Rocca, "An mbse-based requirement verification framework to support the mdao process," in *AIAA Aviation 2022 Forum*, 2022.
- [80] M. Glatt, C. Sinnwell, L. Yi, S. Donohoe, B. Ravani, and J. C. Aurich, "Modeling and implementation of a digital twin of material flows based on physics simulation," *Journal of Manufacturing Systems*, 2021.
- [81] C. Zhang, J. Chen, J. Li, Y. Peng, and Z. Mao, "Large language models for human-robot interaction: A review," *Biomimetic Intelligence and Robotics*, p. 100131, 2023.
- [82] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.
- [83] M. Finlayson, S. Swayamdipta, and X. Ren, "Logits of api-protected llms leak proprietary information," *arXiv preprint arXiv:2403.09539*, 2024.
- [84] N. Yousefnezhad, A. Malhi, and K. Främling, "Security in product lifecycle of iot devices: A survey," *Journal of Network and Computer Applications*, vol. 171, p. 102779, 2020.
- [85] B. Liu, M. Ding, S. Shaham, W. Rahayu, F. Farokhi, and Z. Lin, "When machine learning meets privacy: A survey and outlook," *ACM Computing Surveys (CSUR)*, vol. 54, no. 2, pp. 1–36, 2021.
- [86] M. Rigaki and S. Garcia, "A survey of privacy attacks in machine learning," *ACM Computing Surveys*, vol. 56, no. 4, pp. 1–34, 2023.
- [87] S. Yeom, I. Giacomelli, M. Fredrikson, and S. Jha, "Privacy risk in machine learning: Analyzing the connection to overfitting," in *2018 IEEE 31st computer security foundations symposium (CSF)*. IEEE, 2018, pp. 268–282.
- [88] S. Yeom, I. Giacomelli, A. Menaged, M. Fredrikson, and S. Jha, "Overfitting, robustness, and malicious algorithms: A study of potential causes of privacy risk in machine learning," *Journal of Computer Security*, vol. 28, no. 1, pp. 35–70, 2020.
- [89] O. Podoliaka, V. Mushkatblat, and A. Kaplan, "Privacy attacks based on correlation of dataset identifiers: Assessing the risk," in *2022 IEEE 12th Annual Computing and Communication Workshop and Conference (CCWC)*. IEEE, 2022, pp. 0808–0815.
- [90] H. Liu, Y. Wu, Z. Yu, and N. Zhang, "Please tell me more: Privacy impact of explainability through the lens of membership inference attack," in *2024 IEEE Symposium on Security and Privacy (SP)*. IEEE Computer Society, 2024, pp. 120–120.
- [91] R. Shokri, M. Strobel, and Y. Zick, "On the privacy risks of model explanations," in *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, 2021, pp. 231–241.
- [92] X. Zhao, W. Zhang, X. Xiao, and B. Lim, "Exploiting explanations for model inversion attacks," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 682–692.
- [93] D. Chen, N. Yu, Y. Zhang, and M. Fritz, "Gan-leaks: A taxonomy of membership inference attacks against generative models," in *Proceedings of the 2020 ACM SIGSAC conference on computer and communications security*, 2020, pp. 343–362.
- [94] B. Hilprecht, M. Härterich, and D. Bernau, "Monte carlo and reconstruction membership inference attacks against generative models," *Proceedings on Privacy Enhancing Technologies*, 2019.
- [95] J. Yang, H. Jin, R. Tang, X. Han, Q. Feng, H. Jiang, S. Zhong, B. Yin, and X. Hu, "Harnessing the power of llms in practice: A survey on chatgpt and beyond," *ACM Transactions on Knowledge Discovery from Data*, vol. 18, no. 6, pp. 1–32, 2024.
- [96] J. Lu, X. S. Zhang, T. Zhao, X. He, and J. Cheng, "April: Finding the achilles' heel on privacy for vision transformers," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 10051–10060.
- [97] A. Vakanski, "Adversarial machine learning: Defenses against privacy attacks," <https://www.webpages.uidaho.edu/vakanski>, 2024, accessed: 2024-05-01.
- [98] C. Dwork, A. Roth *et al.*, "The algorithmic foundations of differential privacy," *Foundations and Trends® in Theoretical Computer Science*, vol. 9, no. 3–4, pp. 211–407, 2014.
- [99] C. Dwork, A. Smith, T. Steinke, and J. Ullman, "Exposed! a survey of attacks on private data," *Annual Review of Statistics and Its Application*, vol. 4, pp. 61–84, 2017.
- [100] A. Majeed and S. Lee, "Anonymization techniques for privacy preserving data publishing: A comprehensive survey," *IEEE access*, vol. 9, pp. 8512–8545, 2020.
- [101] J. Li, X. Kuang, S. Lin, X. Ma, and Y. Tang, "Privacy preservation for machine learning training and classification based on homomorphic encryption schemes," *Information Sciences*, vol. 526, 2020.
- [102] T. Chen, H. Bao, S. Huang, L. Dong, B. Jiao, D. Jiang, H. Zhou, J. Li, and F. Wei, "The-x: Privacy-preserving transformer inference with homomorphic encryption," *arXiv preprint arXiv:2206.00216*, 2022.
- [103] F. Miresghallah, M. Taram, P. Vepakomma, A. Singh, R. Raskar, and H. Esmaeilzadeh, "Privacy in deep learning: A survey," *arXiv preprint arXiv:2004.12254*, 2020.
- [104] T. Singh, H. Aditya, V. K. Madiseti, and A. Bahga, "Whispered tuning: Data privacy preservation in fine-tuning llms through differential privacy," *Journal of Software Engineering and Applications*, vol. 17, no. 1, pp. 1–22, 2024.
- [105] F. Boenisch, "A systematic review on model watermarking for neural networks," *Frontiers in big Data*, vol. 4, p. 729663, 2021.
- [106] J. Zhang, D. Chen, J. Liao, W. Zhang, H. Feng, G. Hua, and N. Yu, "Deep model intellectual property protection via deep watermarking," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 8, pp. 4005–4020, 2021.
- [107] S. Neupane, J. Ables, W. Anderson, S. Mittal, S. Rahimi, I. Banicescu, and M. Seale, "Explainable intrusion detection systems (x-ids): A survey of current methods, challenges, and opportunities," *IEEE Access*, vol. 10, pp. 112392–112415, 2022.
- [108] K. Kobayashi and S. B. Alam, "Explainable, interpretable, and trustworthy ai for an intelligent digital twin: A case study on remaining useful life," *Engineering Applications of Artificial Intelligence*, vol. 129, p. 107620, 2024.
- [109] B. Kuipers, "Perspectives on ethics of ai," in *The Oxford Handbook of Ethics of AI*. Oxford University Press, 2020, p. 421.
- [110] I. Asimov, "Three laws of robotics," *Asimov, I. Runaround*, 1941.
- [111] G. Veruggio, "The birth of roboethics," 2005.
- [112] V. Dignum, *Responsible artificial intelligence: how to develop and use AI in a responsible way*. Springer, 2019, vol. 1.
- [113] M. Janssen, P. Brous, E. Estevez, L. S. Barbosa, and T. Janowski, "Data governance: Organizing data for trustworthy artificial intelligence," *Government information quarterly*, vol. 37, no. 3, p. 101493, 2020.
- [114] E. Eryurek, U. Gilad, V. Lakshmanan, A. Kibung'uchi-Grant, and J. Ashdown, *Data governance: The definitive guide*. "O'Reilly", 2021.
- [115] H. Li, L. Yu, and W. He, "The impact of gdpr on global technology development," pp. 1–6, 2019.
- [116] E. Goldman, "An introduction to the california consumer privacy act (ccpa)," *Santa Clara Univ. Legal Studies Research Paper*, 2020.
- [117] M. Mora-Cantalops, S. Sánchez-Alonso, E. García-Barriocanal, and M.-A. Sicilia, "Traceability for trustworthy ai: A review of models and tools," *Big Data and Cognitive Computing*, vol. 5, no. 2, 2021.
- [118] M. Holler, F. Uebernickel, and W. Brenner, "Digital twin concepts in manufacturing industries-a literature review and avenues for further research," in *Proceedings of the 18th international conference on industrial engineering*, 2016.
- [119] G. White, A. Zink, L. Codecá, and S. Clarke, "A digital twin smart city for citizen feedback," *Cities*, vol. 110, p. 103064, 2021.
- [120] L. Chang, L. Zhang, C. Fu, and Y.-W. Chen, "Transparent digital twin for output control using belief rule base," *IEEE Transactions on Cybernetics*, vol. 52, no. 10, pp. 10364–10378, 2022.
- [121] UNSW, "Trusted autonomy," <https://www.unsw.edu.au/canberra/our-research/research-excellence/artificial-intelligence/trusted-autonomy>, 2024, accessed: 2024-05-01.
- [122] I. Pérez-Hurtado, M. Á. Martínez-del Amor, G. Zhang, F. Neri, and M. J. Pérez-Jiménez, "A membrane parallel rapidly-exploring random tree algorithm for robotic motion planning," *Integrated Computer-Aided Engineering*, vol. 27, no. 2, pp. 121–138, 2020.
- [123] W. Jia, M. Sun, J. Lian, and S. Hou, "Feature dimensionality reduction: a review," *Complex & Intelligent Systems*, vol. 8, no. 3, 2022.
- [124] T. Liu, J. Wang, Q. Liu, S. Alibhai, T. Lu, and X. He, "High-ratio lossy compression: Exploring the autoencoder to compress scientific data," *IEEE Transactions on Big Data*, vol. 9, no. 1, pp. 22–36, 2021.
- [125] F. Chazal and B. Michel, "An introduction to topological data analysis: fundamental and practical aspects for data scientists," *Frontiers in artificial intelligence*, vol. 4, p. 108, 2021.