

Cybersecurity Management

T8 – Cybersecurity and AI

2025-2026
Marc Ruiz

marc.ruiz-ramirez@upc.edu



Introduction to *AI/ML*

Basic concepts



ARTIFICIAL INTELLIGENCE

A program that can sense, reason, act, and adapt.



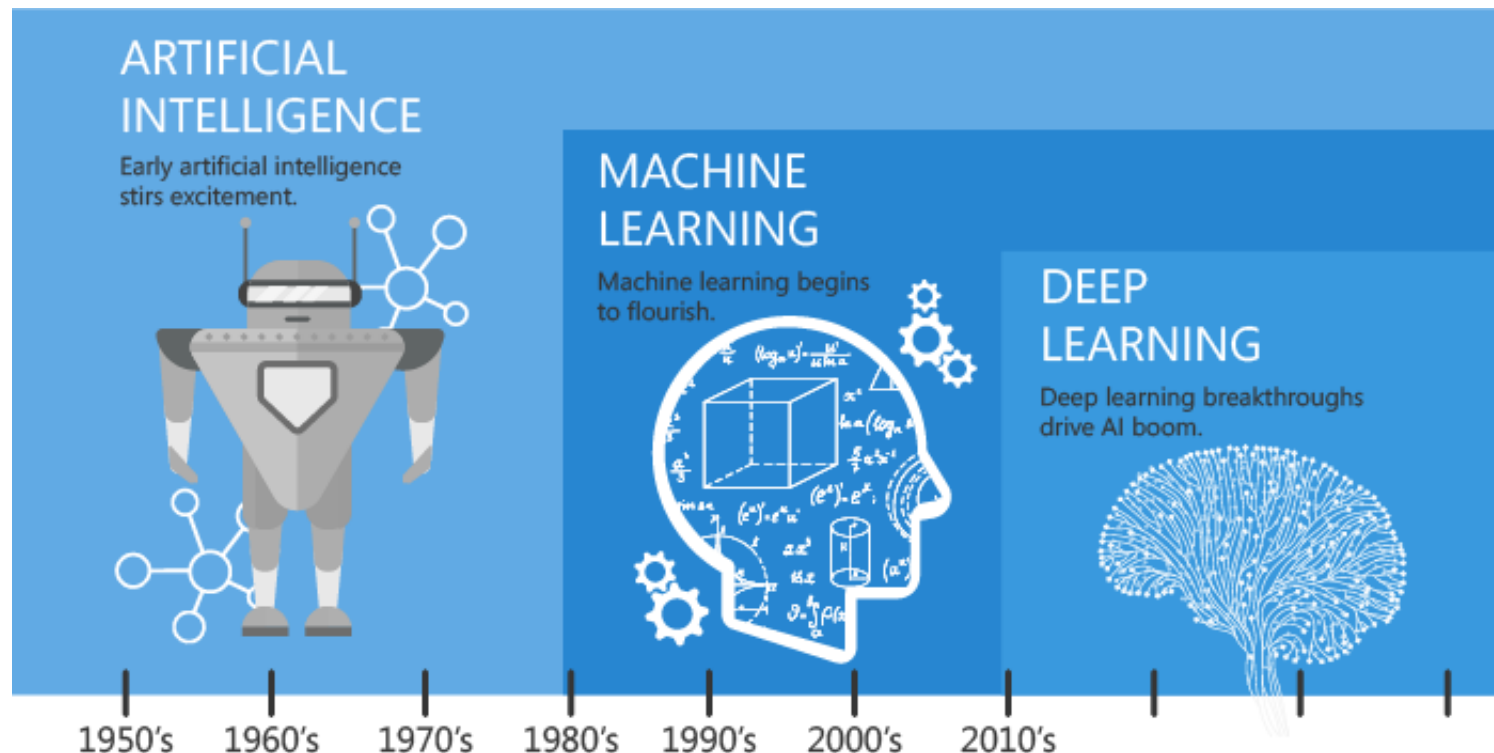
MACHINE LEARNING

Algorithms whose performance improve as they are exposed to more data over time.

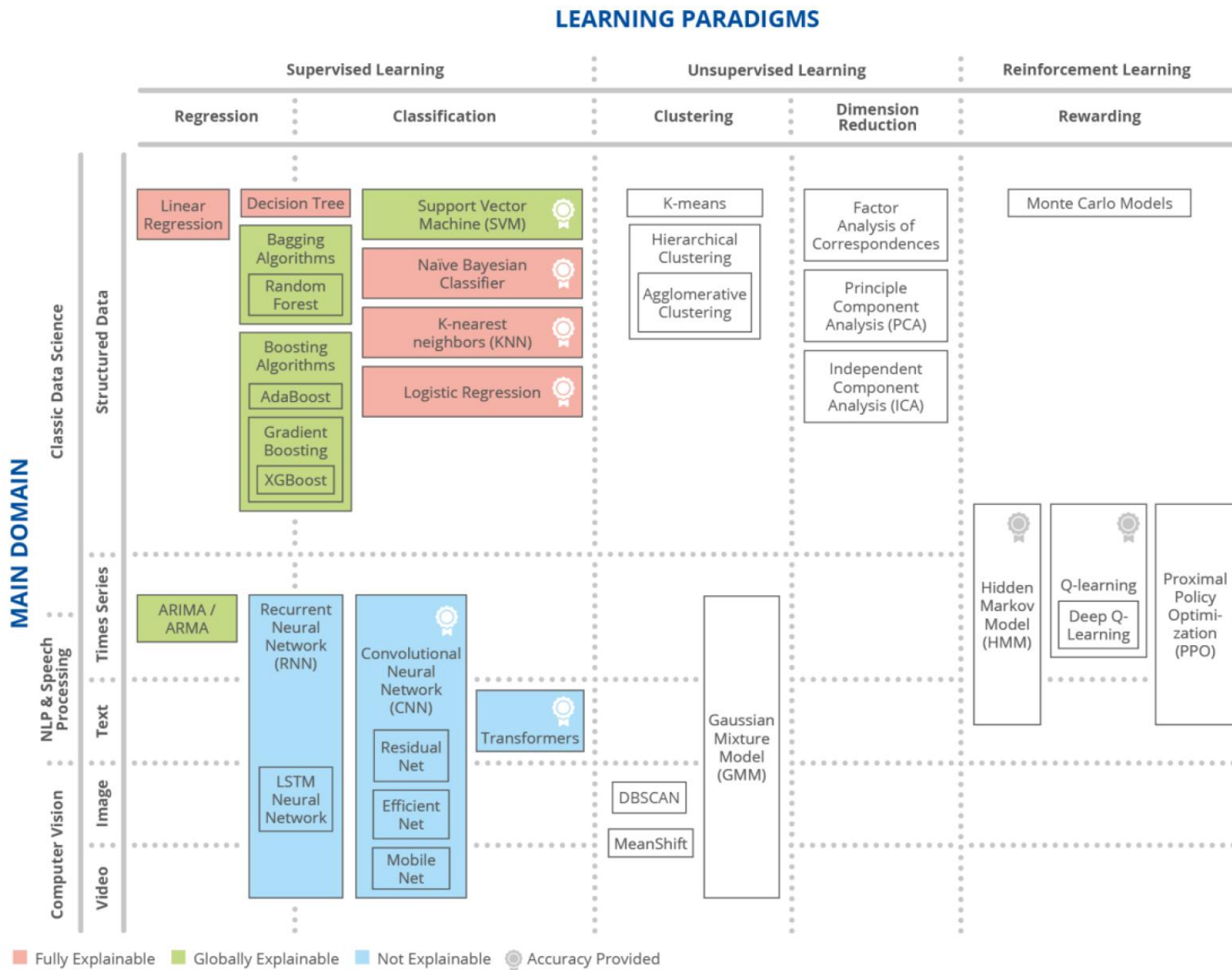


DEEP LEARNING

Subset of machine learning in which multilayered neural networks learn from vast amount of data.



ML taxonomy



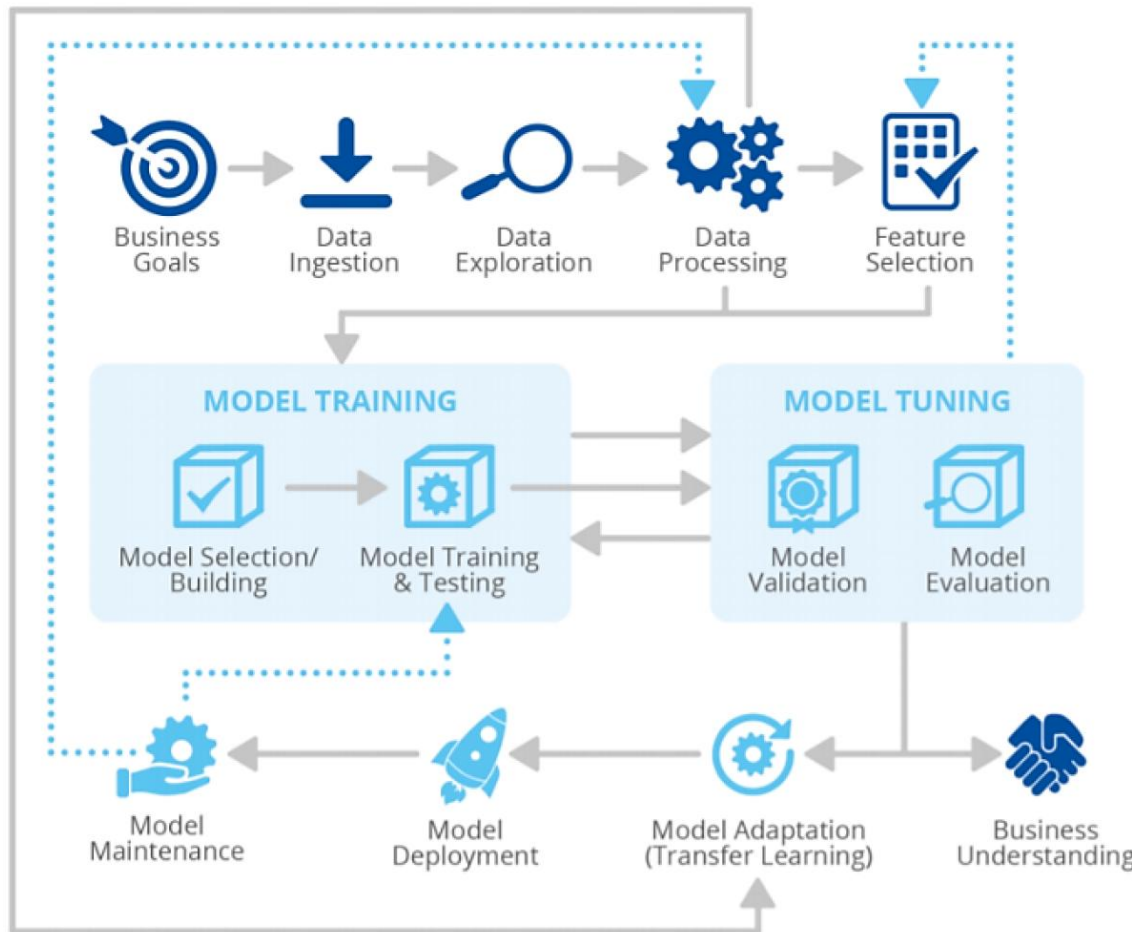
Domains and data types

Main domain	Data type	Definition
Computer Vision	Image	Visual representation of a matrix of pixels constituted of 1 channel for black and white images, 3 elements (RGB) for coloured images or 4 elements (RGBA) for coloured images with opacity.
	Video	A succession of images (frames), sometimes grouped with a time series (a sound).
NLP & Speech processing	Text	A succession of characters (e.g. a tweet, a text field).
	Time series ⁵	A series of data points (e.g. numerical) indexed in time order.
Classic Data Science	Structured Data	<p>Data organised in a predefined model of array with one specific column for each feature (e.g. textual, numerical data, date). To be more accurate, structured data refer to organised data that can be found in a relational data base for example (that may contain textual columns as mentioned).</p> <p>Quantitative data can be distinguished from qualitative data. Quantitative data corresponds to the numerical data that can supports some arithmetic operations whereas qualitative data is usually used as categorical data to classify data according to their similarities.</p>

Learning Paradigms

Learning paradigm	Subtypes	Definition
Supervised learning	Classification	Classification is the process of predicting the class of given data points. (Is the picture a cat or a dog?)
	Regression	Regression models are used to predict a continuous value. (Predict the price of a house based on its features).
Unsupervised learning	Clustering	Clustering is the task of dividing a set of data points into several groups such that data points in the same groups are more similar each other than from the data points of the other groups.
	Dimensionality reduction	Dimensionality reduction refers to techniques for reducing the number of input variables in training data.
Reinforcement learning	Rewarding	Rewarding is an area of ML concerned with how intelligent agents ought to take actions in an environment to maximise the notion of cumulative reward, learning by using feedback from their experiences.

AI/ML lifecycle



DATA COLLECTION

Retrieve data from client's internal storages or external sources



DATA CLEANING

Identify and correct wrong values that may negatively impact an algorithm



DATA PREPROCESSING

Improve data quality by shedding light on relevant information and making it easy to use for ML algorithms: • Dimensionality Reduction • Clustering • Feature Engineering • Data Augmentation • Rescaling



MODEL DESIGN AND IMPLEMENTATION

Choose a predefined model or design a new model and define its parameters



MODEL TRAINING

Train one or a combination of algorithms to accomplish a specific task • Regression • Classification • Clustering • Rewarding



MODEL TESTING

Test the model on unknown data



OPTIMISATION

Apply some technics of hyperparameter tuning to improve the model's performance



MODEL EVALUATION

Define some technical and business metrics to evaluate the model's performance



MODEL DEPLOYMENT

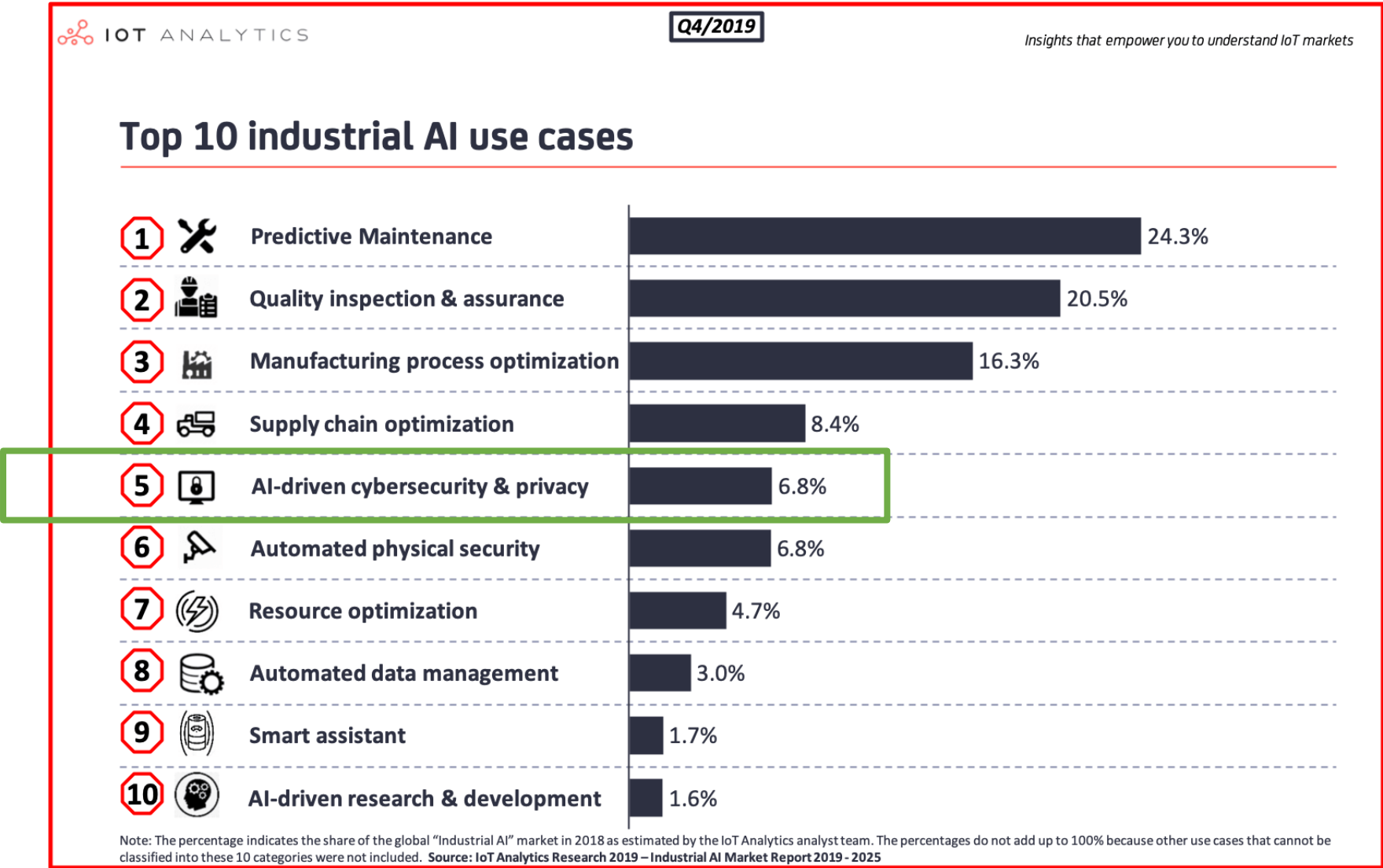
Put the model in production on premise servers or cloud platforms to run and user/model interactions (ex: API)



MONITORING AND INFERENCE

Correspond to the exploitation: observation of the reporting usage of the model and supervision of its performance

AI use cases



Trends



Digital Twins



Represents assets in the physical world with a digital model



Looks and feels like the real environment



Simulates models forward with varying degrees of fidelity

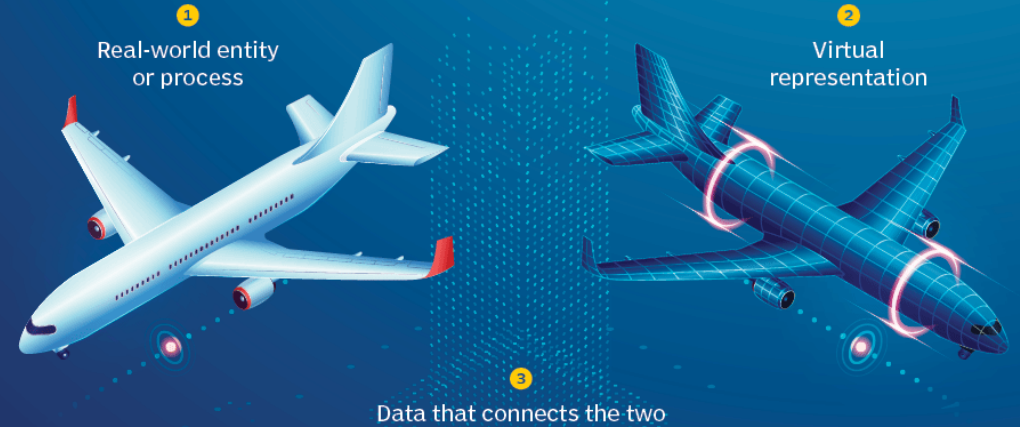


Is NOT just a data model.
It must include relational
interaction



Connects with relevant
time data to ensure the
model mirrors reality

The three elements of a digital twin



Digital Twins

TECHNOLOGIES USED IN DIGITAL TWINS

IoT sensors enable constant data transmission, which is used to create a digital duplicate of the physical object

Due to its visualization capabilities, XR allows to digitally model physical objects

IoT

XR

Cloud

AI

Cloud computing allows to store gained data in the virtual cloud and easily access them from any location

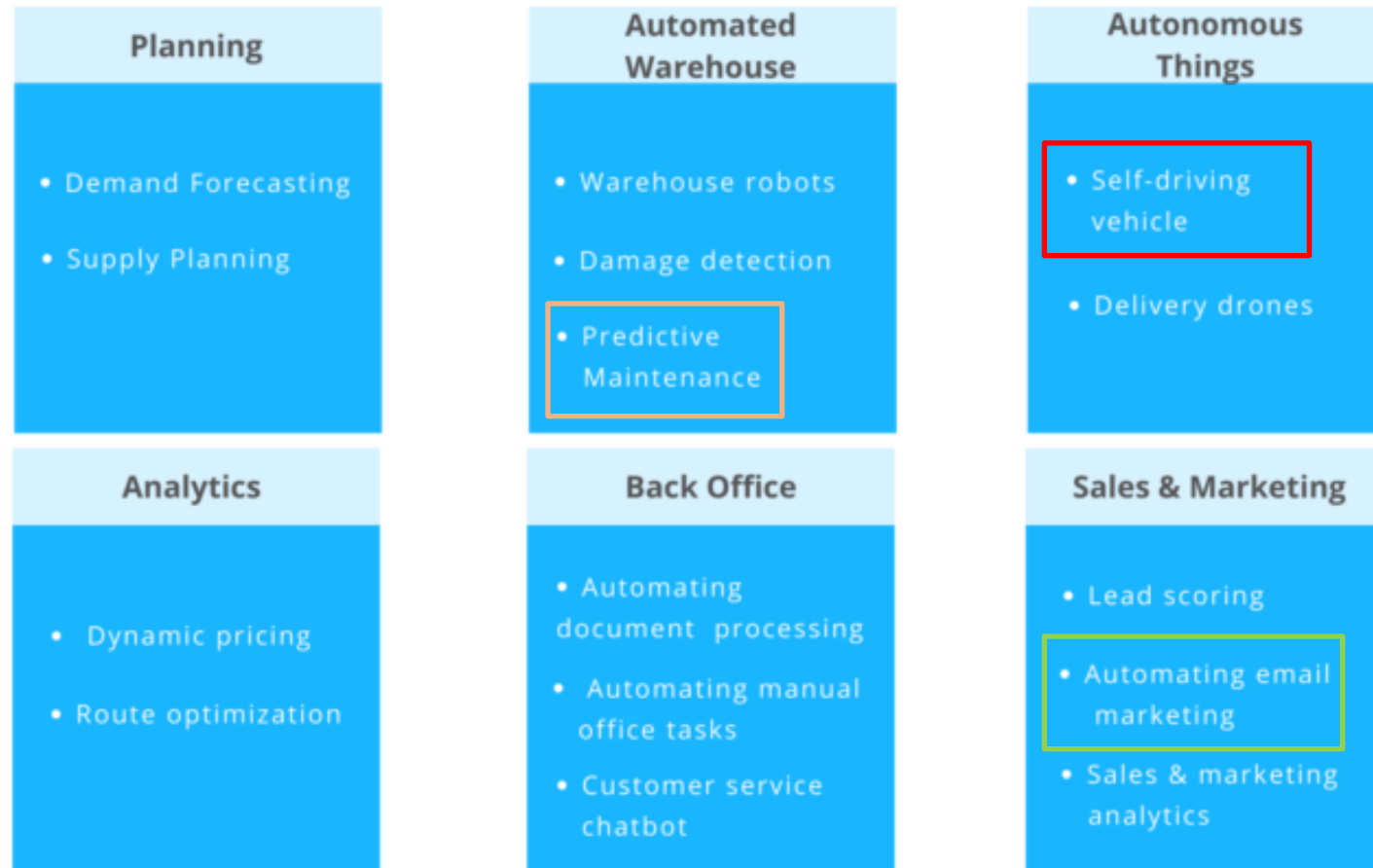
As an advanced analytical tool, AI automatically analyze obtained data, provide valuable insights and make predictions



Applications and risks of AI



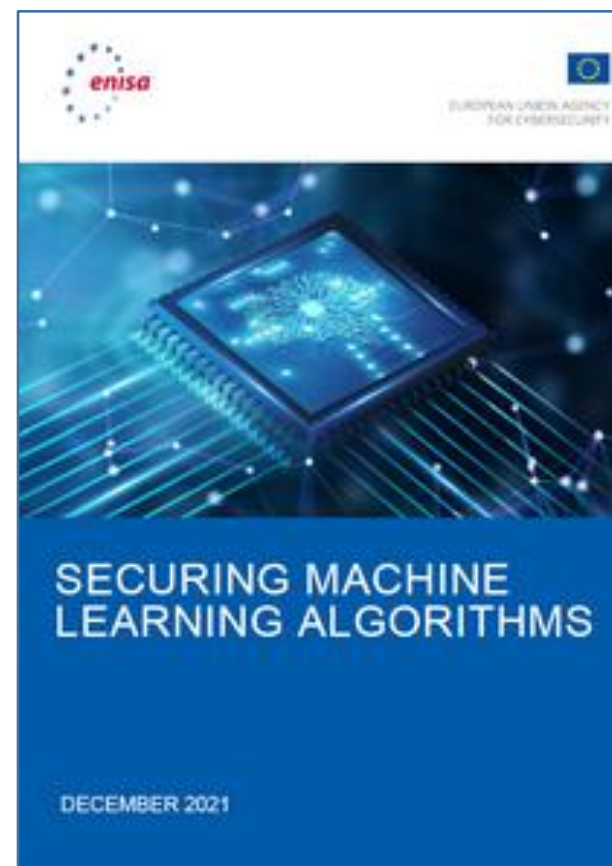
AI Applications in Logistics



Cybersecurity for AI

Analysis of Threats, Vulnerabilities,
and Countermeasures

References



Clarification

- **Cybersecurity for AI**



We focus on this

- All those methods, practices, tools, recommendations, etc related with cybersecurity that can make AI/ML algorithms and procedures more robust and secure
- Example: encrypt a ML model to avoid extracting information from them

- **AI for cybersecurity**

- Using AI-based methods to develop and implement cybersecurity-related tools
- Example: a IDS using artificial neural networks

Selected Threats -> Attacks

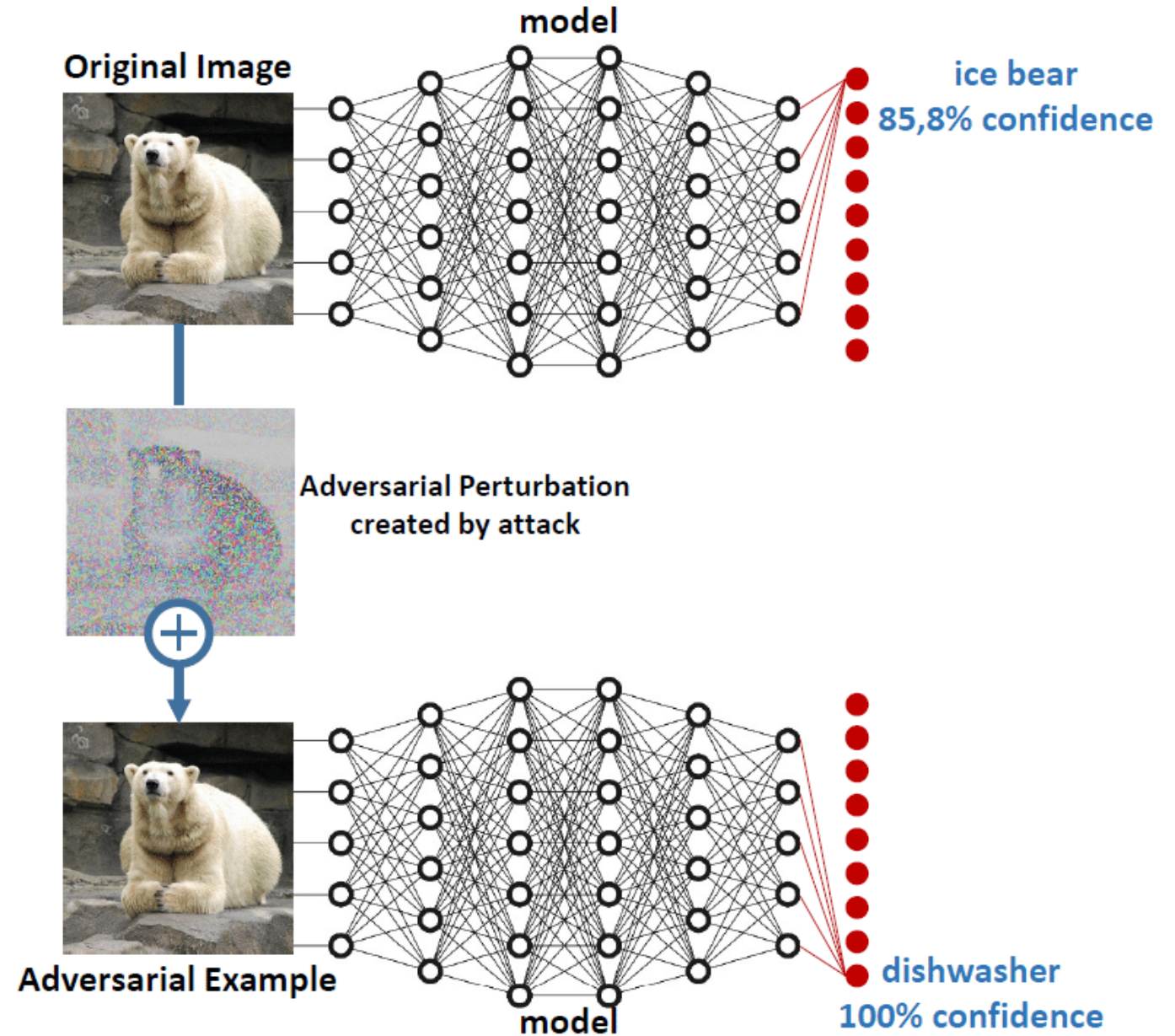
1. Evasion
2. Oracle
3. Poisoning
4. Model/data disclosure
5. Compromise of ML application components
6. Failure or malfunction of ML application

Evasion - Definition



- The attacker works on the ML algorithm's inputs to find small perturbations leading to large modification of its outputs
- a.k.a., adversarial examples.
- Example: the projection of images on a house could lead the algorithm of an autonomous car to take the decision to suddenly make it brake.
- In some cases, the attacker has access to information (model parameters, etc.) that can allow him to directly build adversarial examples.
- Example: use the model's gradient to find the best perturbation to add to the input data to evade the model.

Adversarial Sample



Fezza, Sid Ahmed & Bakhti, Yassine & Hamidouche, Wassim & Déforges, Olivier. (2019). Perceptual Evaluation of Adversarial Attacks for CNN-based Image Classification. 10.1109/QoMEX.2019.8743213.

Evasion - Vulnerabilities



- Lack of detection of abnormal inputs
- Poor consideration of evasion attacks in the model design implementation
- Lack of training based on adversarial attacks
- Using a widely known model allowing the attacker to study it
- Inputs totally controlled by the attacker which allows for input-output-pairs
- Too much information available on the model
- Too much information about the model given in its outputs

Oracle - Definition



- The attacker explores a model by providing a series of carefully crafted inputs and observing outputs.
- Previous steps to more harmful types such as evasion or poisoning
- It is as if the attacker made the model talk to then better compromise it or to obtain information about it
- Example: an attacker studies the set of input-output pairs and uses the results to retrieve training data.

Oracle - Vulnerabilities



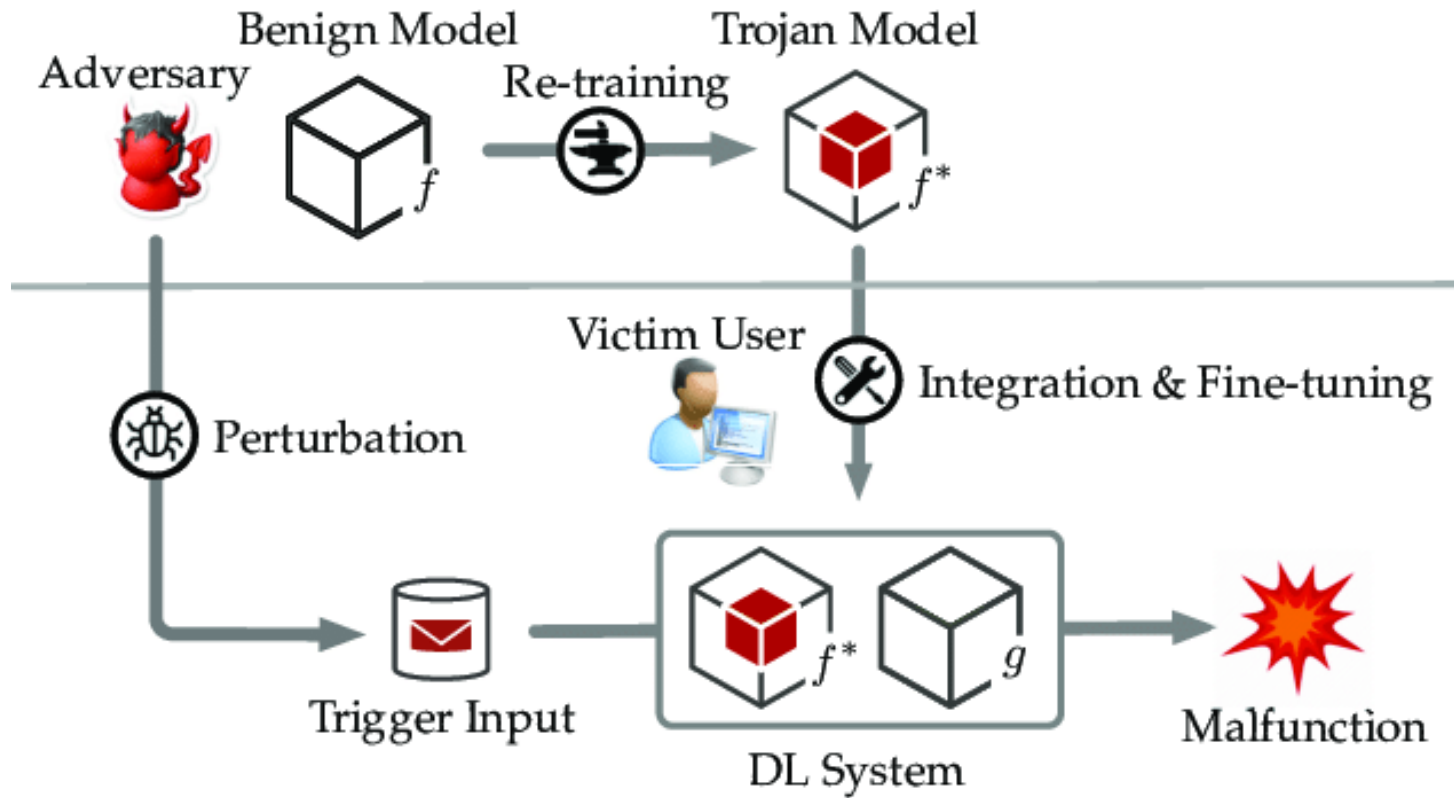
- Poor access rights management
- The model allows private information to be retrieved
- Too much information about the model given in its outputs
- Too much information available on the model
- Lack of consideration of attacks to which ML applications could be exposed to
- Lack of security process to maintain a good security level of the components of the ML application
- Weak access protection mechanisms for ML model components

Poisoning - Definition



- A type of attack in which the attacker altered data or model to modify the ML algorithm's behavior in a chosen direction (backdoor)
- It is as if the attacker conditioned the algorithm according to its motivations.
- Such attacks are also called causative attacks.
- Example: massively indicating to an image recognition algorithm that images of dogs are indeed cats to lead it to interpret it this way.
 - Label modification

Backdoor attack



Pang, Ren & Zhang, Zheng & Gao, Xiangshan & Xi, Zhaohan & Ji, Shouling & Cheng, Peng & Wang, Ting. (2020). TROJANZOO: Everything you ever wanted to know about neural backdoors (but were afraid to ask).

Poisoning - Vulnerabilities



- Model easy to poison
- Lack of data for increasing robustness to poisoning
- Poor access rights management
- Poor data management
- Undefined indicators of proper functioning, making complex compromise identification
- Lack of consideration of attacks to which ML applications could be exposed to
- Use of uncontrolled data
- Use of unsafe data or models (e.g. with transfer learning)
- Lack of control for poisoning
- No detection of poisoned samples in the training dataset
- Weak access protection mechanisms for ML model components

Data disclosure - Definition



- Leak of data manipulated by ML algorithms.
- Reasons of data leakage
 - inadequate access control
 - a handling error of the project team
 - the entity that owns the model and the entity that owns the data are distinct.
- To train the model, sharing the data with model provider (third-party) might lead to sharing sensitive data.

Model disclosure - Definition



- Leak of the internals (i.e. parameter values) of the ML model.
- Reasons
 - human error
 - contraction with a third party with a too low security level.

Model/Data disclosure - Vulnerabilities



- Poor access rights management
- Existence of unidentified disclosure scenarios
- Weak access protection mechanisms for ML model components
- Lack of security process to maintain a good security level of the components of the ML application
- Unprotected sensitive data on test environments
- Too much information about the model given in its outputs
- The model can allow private information to be retrieved
- Disclosure of sensitive data for ML algorithm training
- Too much information available on the model
- Too much information about the model given in its outputs

Compromise of ML application components - Definition



- Compromise of a component or developing tool of the ML application.
 - Example: compromise of one of the open-source libraries used by the developers to implement the ML algorithm.

Compromise of ML application components - Vulnerabilities



- Poor access rights management
- Too much information available on the model
- Existence of several vulnerabilities because the ML application was not included into process for integrating security into projects
- Use of vulnerable components (among the whole supply chain)
- Too much information about the model given in its outputs
- Existence of unidentified compromise scenarios
- Undefined indicators of proper functioning, making complex compromise identification
- Bad practices due to a lack of cybersecurity awareness
- Lack of security process to maintain a good security level of the components of the ML application

Failure or malfunction of ML application - Definition

Stage of the lifecycle

Data Collection	Data Cleaning	Data Preprocessing	Model design	Model Training	Model Testing	Optimisation	Model Evaluation	Model Deployment	Monitoring
-----------------	---------------	--------------------	--------------	----------------	---------------	--------------	------------------	------------------	------------

- This threat refers to ML application failure (e.g. denial of service due to bad input, unavailability due to a handling error).
- The different stakeholders of the model can make mistakes that result in a failure or malfunction of ML application (human error).
 - Example: due to lack of documentation, they may use the application in use-cases not initially foreseen.
- Input data whose format is inappropriate (denial of service).
 - Example: a malicious user of the model constructs an input data (a sponge example) specifically designed to increase the computation time of the model and thus potentially cause a denial of service.

Failure or malfunction of ML application - Vulnerabilities

Stage of the lifecycle

Data Collection	Data Cleaning	Data Preprocessing	Model design	Model Training	Model Testing	Optimisation	Model Evaluation	Model Deployment	Monitoring
-----------------	---------------	--------------------	--------------	----------------	---------------	--------------	------------------	------------------	------------

- Existing biases in the ML model or in the data
- ML application not integrated in the cyber-resilience strategy
- Existence of unidentified failure scenarios
- Undefined indicators of proper functioning, making complex malfunction identification
- Lack of explainability and traceability of decisions taken
- Lack of security process to maintain a good security level of the components of the ML application
- Existence of several vulnerabilities because ML specificities are not integrated in existing policies
- Contract with a low security third party

Security Controls (Countermeasures)

- ENISA identified a list of 37 security controls
- Organisational and Policy
 - More traditional security controls, either organisational or linked to security policies.
- Technical
 - More classic technical security controls.
- **Specific to ML**
 - **Security controls that are specific to applications using ML.**

Organisational

- Apply a RBAC model, respecting the least privileged principle
 - RBAC = Role Based Access Control
- Apply documentation requirements to AI projects
- Assess the regulations and laws the ML application must comply with
- Ensure ML applications comply with data security requirements
- Ensure ML applications comply with identity management, authentication, and access control policies
- Ensure ML applications comply with security and protection policies and are integrated to security operations processes
- Include ML applications into detection and response to security incident processes
- Include ML applications in asset management processes

Technical

- Assess the exposure level of the model used
- Check the vulnerabilities of the components used so that they have an appropriate security level
- Conduct a risk analysis of the ML application
- Control all data used by the ML model
- Define and monitor indicators for proper functioning of the model
- Ensure appropriate protection is deployed for test environments
- Ensure ML applications comply with third parties' security requirements
- Ensure ML projects follow the global process for integrating security into projects

Specific to ML (1/2)

- Add some adversarial examples to the training dataset
- Apply modifications on inputs
- Build explainable models
- Choose and define a more resilient model design
- Enlarge the training dataset
- Ensure that models are unbiased
- Ensure that models respect differential privacy to a sufficient degree
- Ensure that the model is sufficiently resilient to the environment in which it will operate.
- Implement processes to maintain security levels of ML components over time

Specific to ML (2/2)

- Implement tools to detect if a data point is an adversarial example or not
- Integrate ML specificities to awareness strategy and ensure all ML stakeholders are receiving it
- Integrate poisoning control after the "model evaluation" phase
- Use less easily transferable models
- Reduce the available information about the model
- Reduce the information given by the model
- Use federated learning to minimize risk of data breaches
- Use less easily transferable models

Lessons learned...so far

- There is no silver bullet for mitigating ML-specific attacks
 - Some security controls may be bypassed by adaptive attackers.
 - Applied mitigations can still raise the bar for attackers.
- ML-specific mitigation controls are not generally evaluated in a standardised way even if it is a current and important issue to enable comparability.
 - More research should be devoted to standardised benchmarks for comparing ML-specific mitigations on a level playing field.
 - These benchmarks should also be enforced to ensure that the methods used in practice are the ones that perform best.
- Deploying security controls often leads to a trade-off between security and performance
 - This is a topic of particular importance that should be further pursued by the research and cybersecurity communities.

Generative AI (GenAI) and Cybersecurity

GenAI Security Topics

- ETL 2024 -> Emerging topic
 - <https://www.enisa.europa.eu/publications/enisa-threat-landscape-2024>
 - Section 7.6 AI And The Surge Of AI Chatbots: Source And Target Of Attacks
- AI chatbots (OpenAI ChatGPT, Microsoft Copilot and Google Bard) are powerful tools in the hands of cyber-attackers aiming at data breaches...
- ... and at the same time, a powerful tool in the hands of cyber defenders (GenAI for security)...
- ... and also a preferred target for cybercriminals as they are very susceptible to prompt injection and data poisoning, such as malicious data injection into the training datasets
- The EU has put forward a proposal to regulate AI called AI Act
 - <https://digital-strategy.ec.europa.eu/en/policies/european-approach-artificial-intelligence>

GenAI Security Topics

- During the Immersive Labs Prompt Injection Challenge in June-September 2023, 88% of prompt injection challenge participants successfully tricked the GenAI bot into giving away sensitive information.
- This still happens when chatbots are strengthened with security measures, with users crafting more complex prompts forcing GenAI into revealing confidential information
 - No protocols exist to fully prevent prompt injection attacks, nor a clear definition of liability when generative AI is used
- Source code is the most common type of sensitive data shared in ChatGPT
 - 158 incidents per 10,000 users a month
- Followed by regulated data (e.g. financial and healthcare data, PII), intellectual property excluding source code, and passwords and keys, usually embedded in source code.

Malicious LLMs

- <https://www.trustwave.com/en-us/resources/blogs/spiderlabs-blog/wormgpt-and-fraudgpt-the-rise-of-malicious-llms/>
- WormGPT (2021)
 - The hacker chatbot is devoid of any restrictions preventing it from answering questions about illegal activity, unlike mainstream LLMs like ChatGPT.
 - The relatively outdated open source large GPT-J language model from 2021 was used as a platform for creating the chatbot.
 - The chatbot was trained in materials related to malware development, which is how WormGPT was born.
 - E.g., WormGPT writes malware on Python according to malicious requirements, or write a phishing email according to the requirements
 - The seller offers a newer version of WormGPT, WormGPT v2, for €550 Euros annually, and a private build for €5000 Euros which includes access to WormGPT v2.

Malicious LLMs

- <https://www.trustwave.com/en-us/resources/blogs/spiderlabs-blog/wormgpt-and-fraudgpt-the-rise-of-malicious-llms/>
- FraudGPT (2023)
 - Advertised on several Dark Web boards and Telegram channels
 - The price range starts from \$90 – \$200 USD for a monthly subscription, 3 months for \$230 – \$450 USD, \$500 – \$1,000 USD for half a year's subscription, and \$800 – \$1,700 USD for a yearly subscription.
 - The author described FraudGPT as a great tool for creating undetectable malware, writing malicious code, finding leaks and vulnerabilities, creating phishing pages, and for learning hacking

Malicious LLMs

- <https://www.trustwave.com/en-us/resources/blogs/spiderlabs-blog/wormgpt-and-fraudgpt-the-rise-of-malicious-llms/>
- There is not much difference between WormGPT and ChatGPT (with proper build requests) to achieve the desired results.
- In other words, “playing conveniently” with regular LLM engines can lead to malicious purposes
- This is a clear sign that GenAI can become a weapon in the hands of cybercriminals

Cybersecurity Management

T8 – Cybersecurity and AI

2025-2026
Marc Ruiz

marc.ruiz-ramirez@upc.edu

