**Two general categories of risk:**

**Inherent:** e.g., unwanted bias, hallucinations, errors in the generated data, implementation flaws in the model, cybersecurity flaws in the platform on which the AI/ML models is deployed. Dealt with in other standards, e.g.,

1. NIST SP 1270 "Towards a Standard for Identifying and Managing Bias in Artificial Intelligence".
2. NIST SSDF Companion for LLMs – coming soon.



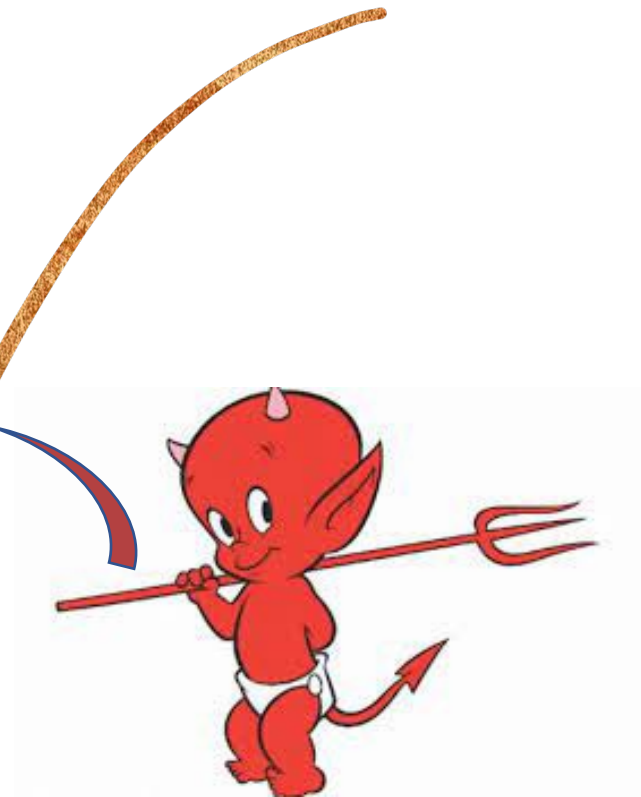**Toxicity in AI Text Generation**

Hi, how are you?

Human

#!?*@

Chatbot

What just happened?

Human

Graphic credit: Julia Nikulski, Towards Data Science

**Adversarial:** deliberate actions by motivated experienced adversaries aiming to

**disrupt, evade, compromise, or abuse**

the operation of the model or its output.

❖ A taxonomy of attacks and mitigations

**A new standard NIST AI 100-2e2023**

**Maintained annually**
- *NIST AI 100-2e2024 ipd – to appear mid-2024*
- *NIST AI 100-2e2024*
- *etc.*

**NIST will seek comments and recommendations on:**
- *What are the latest attacks on the existing AI models?*
- *What are the latest mitigations?*
- *What are the latest trends in AI technologies that promise to transform the industry/society? What potential vulnerabilities do they come with? What promising mitigations may be developed for them?*
- *Is there new terminology that needs standardization?*

NIST Trustworthy and Responsible AI
NIST AI 100-2e2023

**Adversarial Machine Learning**
*A Taxonomy and Terminology of Attacks and Mitigations*

Apostol Vassilev
*Computer Security Division*
*Information Technology Laboratory*

Alina Oprea
*Northeastern University*

Alie Fordyce
Hyrum Anderson
*Robust Intelligence, Inc.*

This publication is available free of charge from:
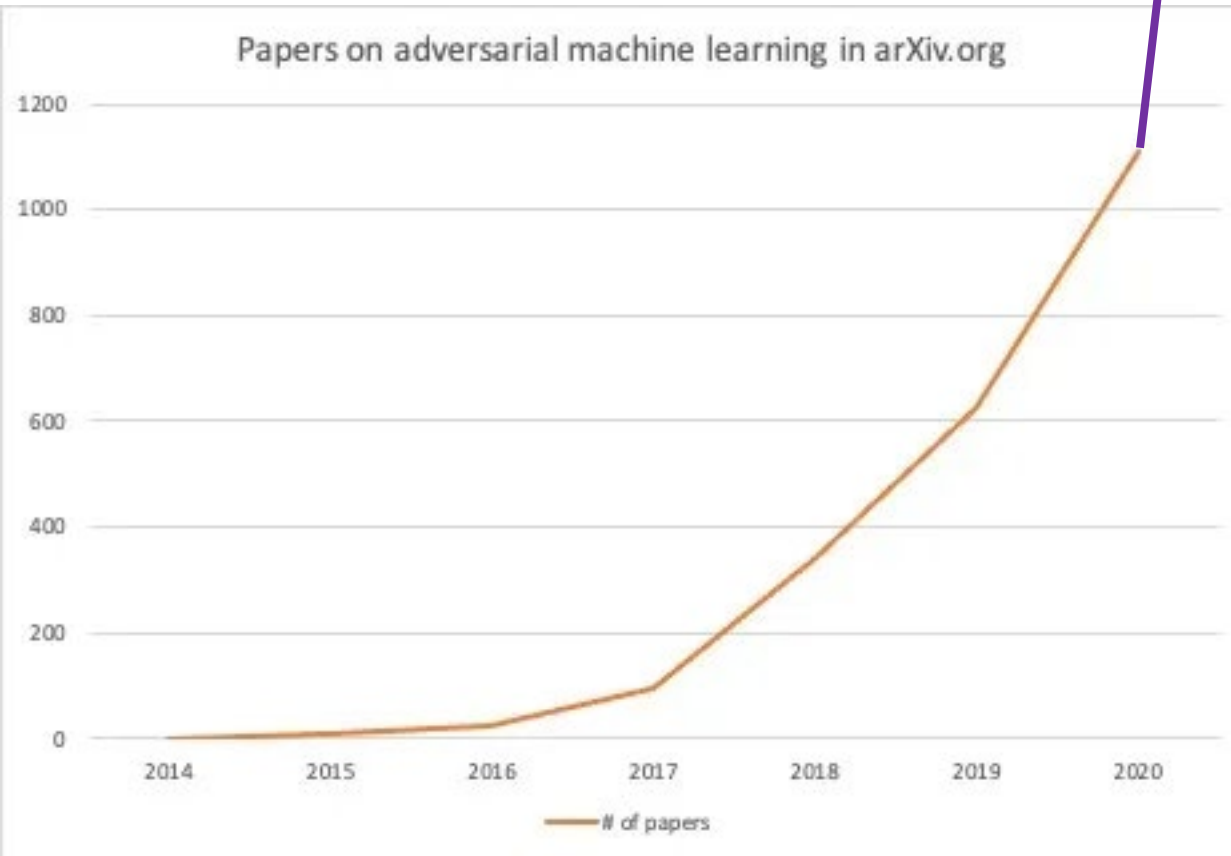https://doi.org/10.6028/NIST.AI.100-2e2023

January 2024

U.S. Department of Commerce
*Gina M. Raimondo, Secretary*

National Institute of Standards and Technology
*Laurie E. Locascio, NIST Director and Under Secretary of Commerce for Standards and Technology*

# AML Pace

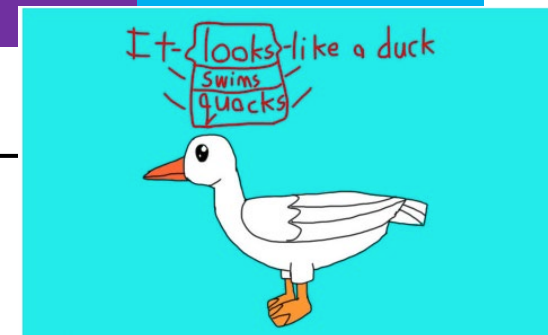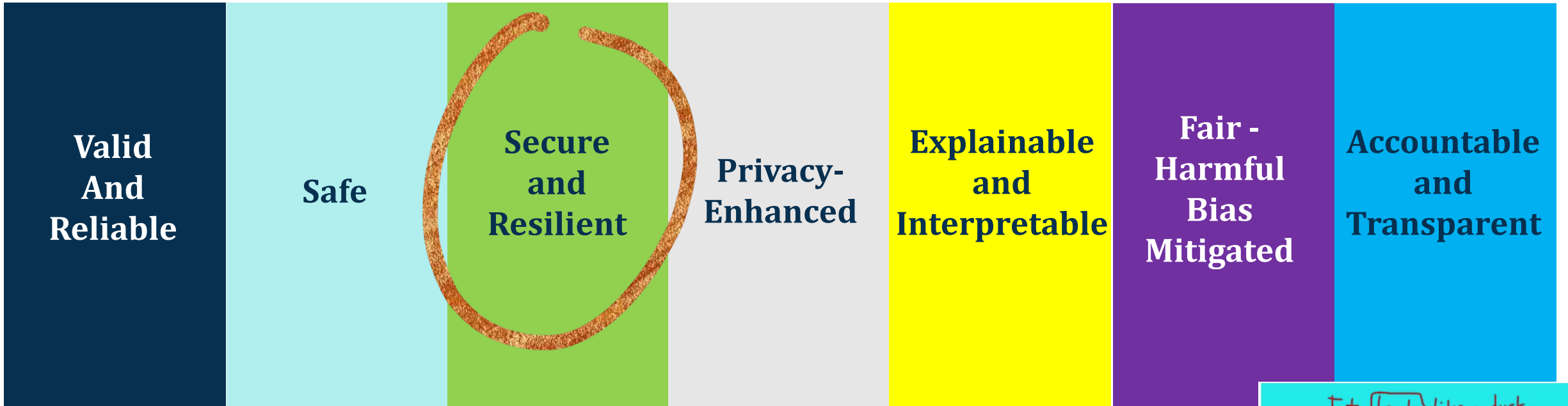Papers on adversarial machine learning in arXiv.org



A search on arXiv for AML articles in
**2021** and **2022**
yielded more than **5,000** references

What drives this enormous growth?

No *information-theoretic* security guarantees for AI algorithms !

Worse, information-theoretic ***impossibility*** results have been established, making the security problem intractable in the existing AI paradigm.

Credit: Ben Dickson https://www.kdnuggets.com/2021/01/machine-learning-adversarial-attacks.html

# Trustworthy AI

❖ The Seven Attributes of Trustworthiness

| Valid And Reliable | Safe | Secure and Resilient | Privacy-Enhanced | Explainable and Interpretable | Fair - Harmful Bias Mitigated | Accountable and Transparent |
|---|---|---|---|---|---|---|

It looks like a duck
swims
quacks

# Trustworthy AI Attributes

❖ Relationships between Attributes

*Accuracy, Fairness, Explainabilty: How do they relate to **Privacy** and **Adversarial Robustness**?*

❖ It is <u>not possible</u> to simultaneously maximize the performance of the AI system with respect to these attributes.

   ❖ **Accuracy** vs. **Adversarial Robustness** tradeoff (          )

   ❖ **Fairness** vs. **Adversarial Robustness**  (          )

   ❖ **Explainability** vs. **Adversarial Robustness**  (          )

   ❖ **Privacy** vs **Fairness** (          )

❖ Organizations need to <u>accept trade-offs</u> and decide priorities depending on the AI system, the use-case, economic, environmental, social, cultural, political, and global implications of the AI technology.
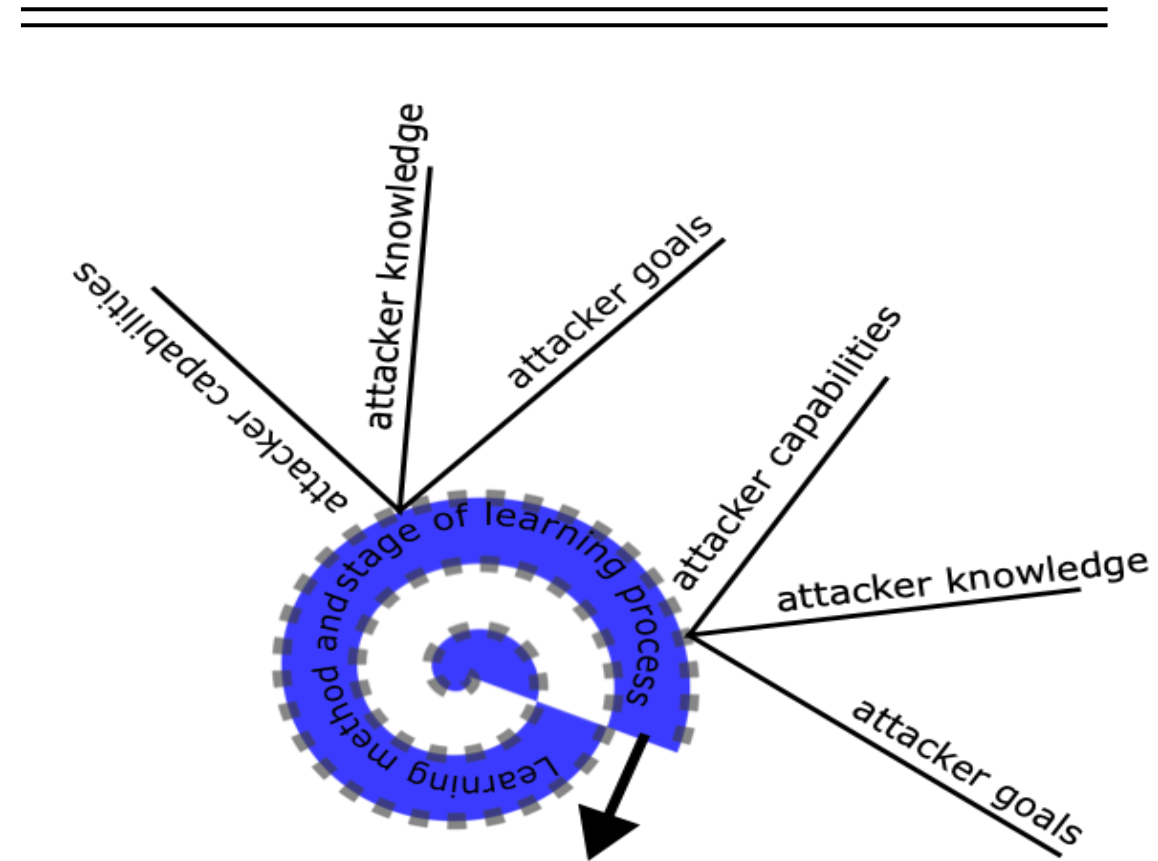
NIST

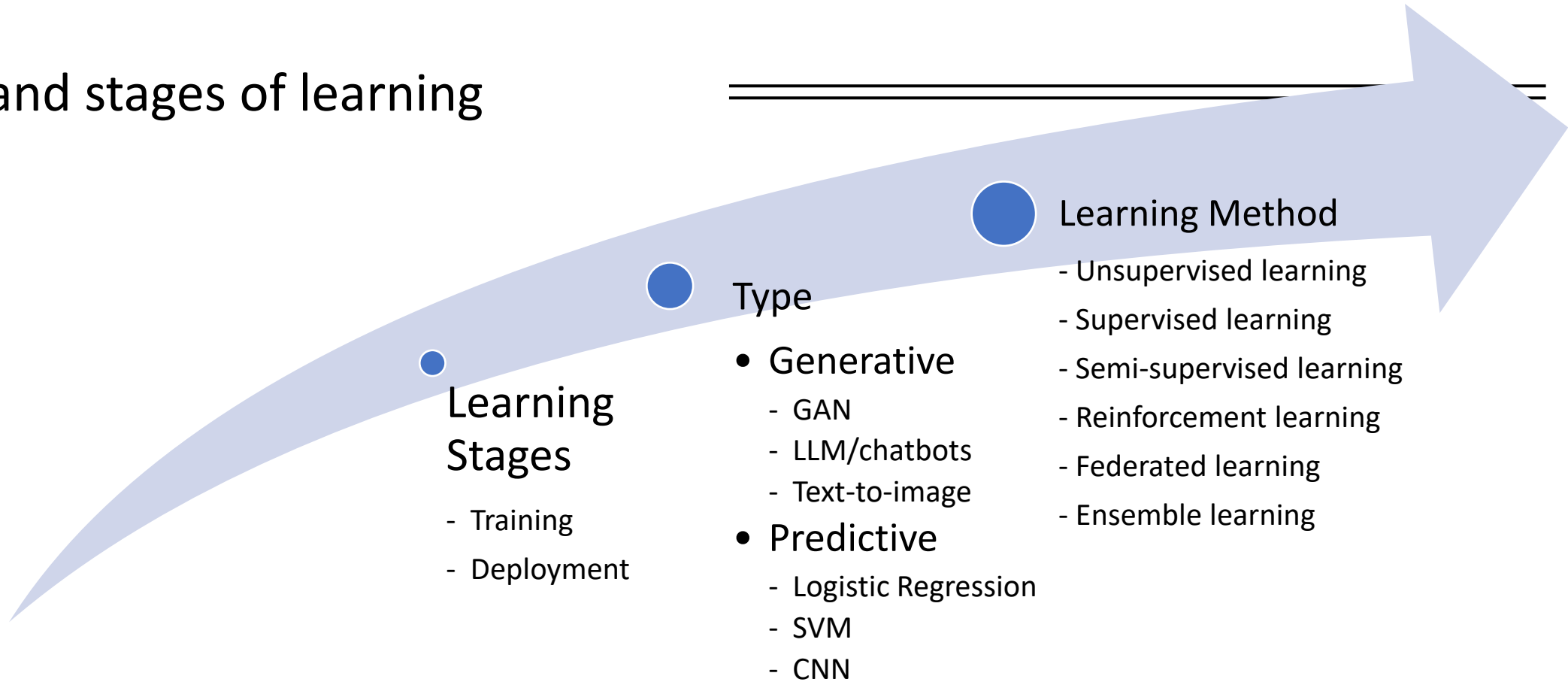❖ A taxonomy of attacks and mitigations

**Four dimensions:**

❖ *Learning method and stage of learning process*

❖ *Attacker goals/objectives*

❖ *Attacker capabilities*

❖ *Attacker knowledge*

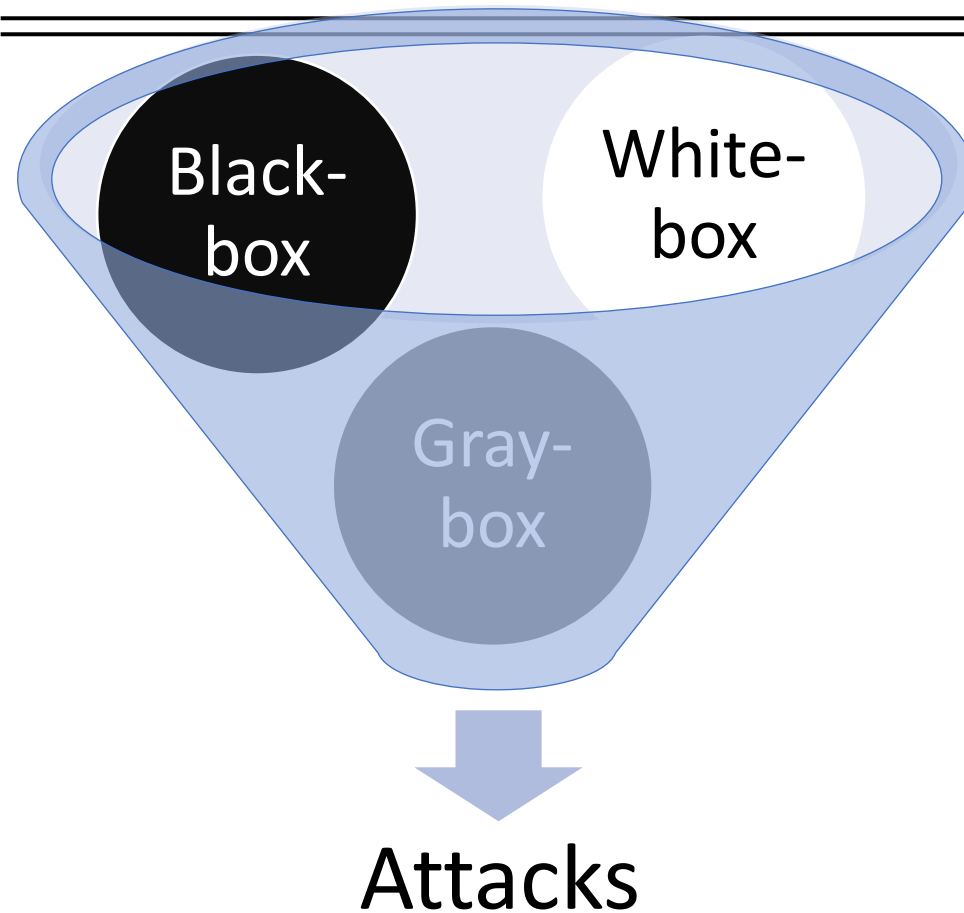**ML models can be attacked at all stages of their lifecycle**

❖ *from design to training to deployment and use*

# Adversarial ML (AML)

❖ Methods and stages of learning

**Learning Method**
- Unsupervised learning
- Supervised learning
- Semi-supervised learning
- Reinforcement learning
- Federated learning
- Ensemble learning

Type
- Generative
  - GAN
  - LLM/chatbots
  - Text-to-image
- Predictive
  - Logistic Regression
  - SVM
  - CNN

Learning Stages
- Training
- Deployment

# Adversarial ML (AML)

❖ Attacker knowledge



Black-box

White-box

Gray-box

Attacks

# Adversarial ML (AML)
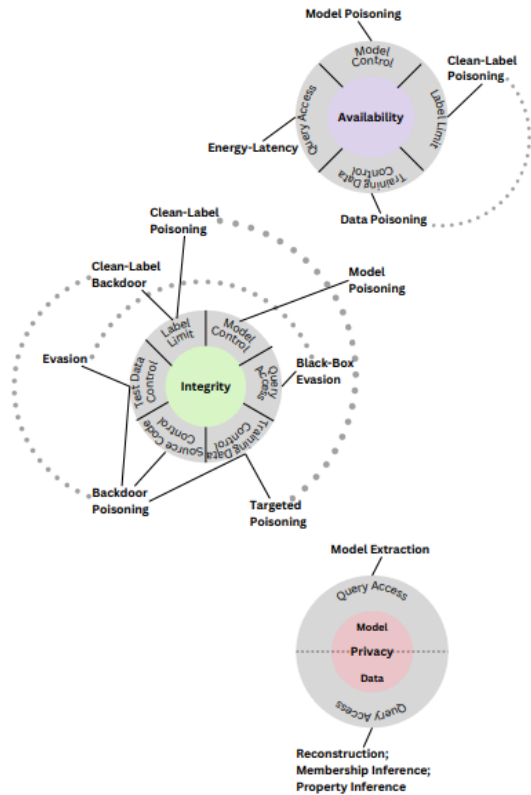
❖ Attacker goals/objectives perspective



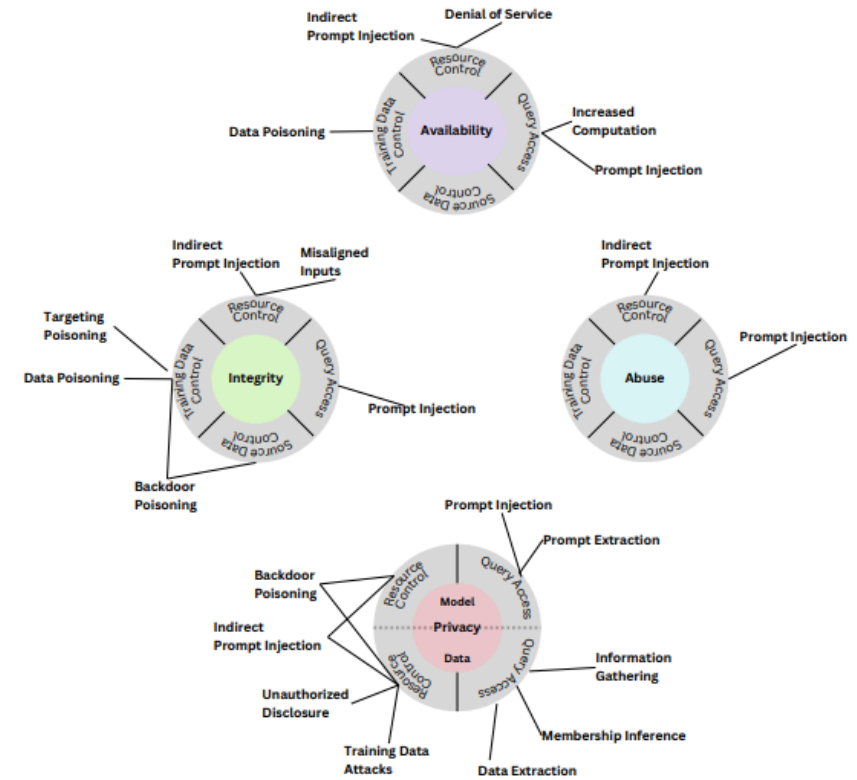Figure 1. Taxonomy of attacks on Predictive AI systems.

Figure 2. Taxonomy of attacks on Generative AI systems

❖ Physical Evasion attack example

**Credit:** Jing et al., "Too Good to Be Safe: Tricking Lane Detection in Autonomous Driving with Crafted Perturbations", USENIX 2021.

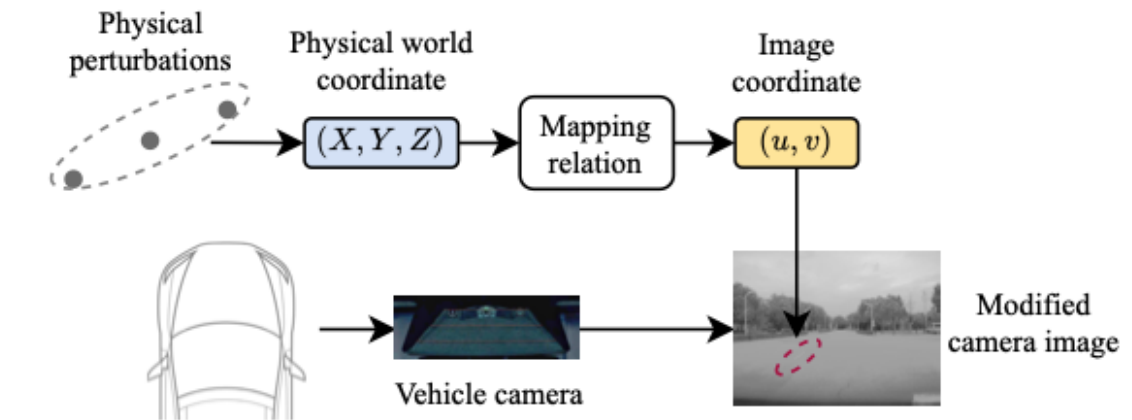**Human Eye Invisibale/Neglectible markings on road cause the vehicle the veer off into the opposite traffic lane**
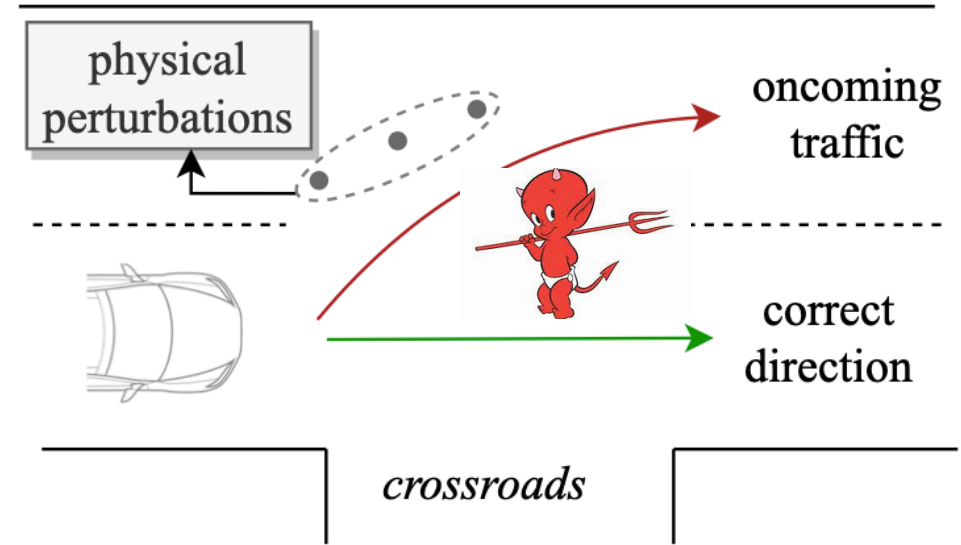


Figure 4: Mapping the coordinate of $(X, Y, Z)$ on markings in physical world to the coordinate of $(u, v)$ on perturbations in digital world.
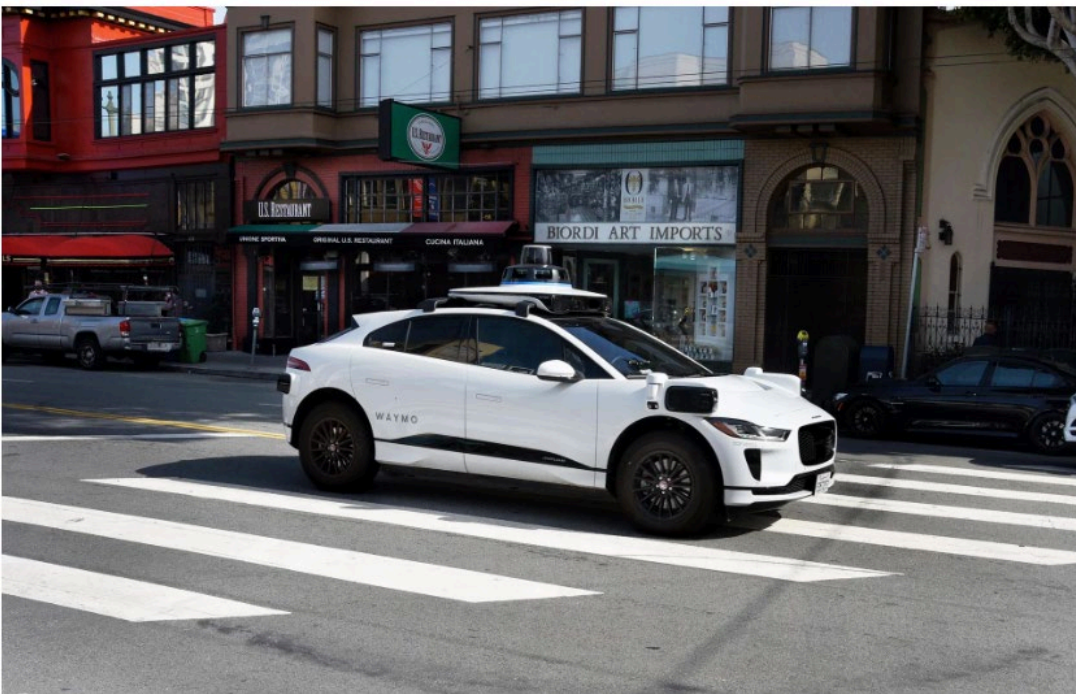
# PredAI AML – the risks are not just anecdotal

NIST

❖ It is not just one company



Autopilot crash, Walnut Creek, CA, 02/18/2023

**NHTSA report June 2023:** Autonomous Driving Systems safety record is currently **lagging** human driver performance for the same number of traveled miles:

**Tesla Autopilot:**
**736 Crashes since 2019,**
**17 of them were fatal and**
**11 deaths have occurred since May 2022**



NHTSA    Ratings    Recalls    Risky Driving

← LAWS & REGULATIONS

**Standing General Order on Crash Reporting**

For incidents involving ADS and Level 2 ADAS

## ❖ Adversarial Training (AT)

- ❖ The most robust approach known so far
- ❖ Due to Goodfellow et al. in 2015
- ❖ Improved by Madry et al. in 2018



## ❖ **But,**

- ❖ In automotive setting AT is reactive by construction:
  - ❖ not all road/traffic conditions leading to incidents are known in advance.
- ❖ actual accident data is fed into the training of the next AI model

Cognitive task automation

≠

cognitive intelligence

**For further info:** see the NIST Automated Vehicle Program

Image credit: Zhao et al., "Adversarial Training Methods for Deep Learning: A Systematic Review, MDPI, 2022.

❖ **Certifiable Robustness**

**Definition:** A classifier is said to be _certifiably robust_ if for any input **x**, one can guarantee that the classifier's prediction is constant _within some set_ around **x**, often an $L_2$ or $L_\infty$ ball.

- In the context of **Lp** norm-bounded perturbations, for a classifier **g**, input **x**, and radius **r**,

$$g(x) = g(x + \delta), \text{ for any perturbation } \delta \text{ such that } \delta \leq r.$$

Given an input (e.g., image  **x** correctly classified by a neural network an adversary can engineer an adversarial perturbation **ε** so small that **x + ε** looks just like **x** to humans, yet
$$g(x) \neq g(x + \varepsilon) \text{ - an incorrect class.}$$

- - the relationship between **ε** and **r** is not absolute – what is invisible to the human eye ( 👁 ) may still be too big for AI{ 🧠 }



"panda"
57.7% confidence

$+\epsilon$

$=$

"gibbon"
99.3% confidence

# Chatbots

**NIST**

❖ Training pipeline



| Pre- Training | Supervised Fine-Tunning | Reward Modeling | Reinforcement Learning |
|---|---|---|---|
| **Raw internet data** low quality/large quantity | **Sets of ideal labeled assistant responses** low quantity/high quality | **Comparisons** human-written, low quantity/high quality | **Human prompts** low quantity/high quality |
| **Language modeling:** predict the next token | **Language modeling:** predict the next token | **Binary classification** predict rewards according to preferences | **Reinforcement learning** generate tokens that maximize the reward |
| Base Model | SFT Model | RM Model (not released) | RLHF Model |

# Chatbots in the enterprise

❖ LLM project pipeline

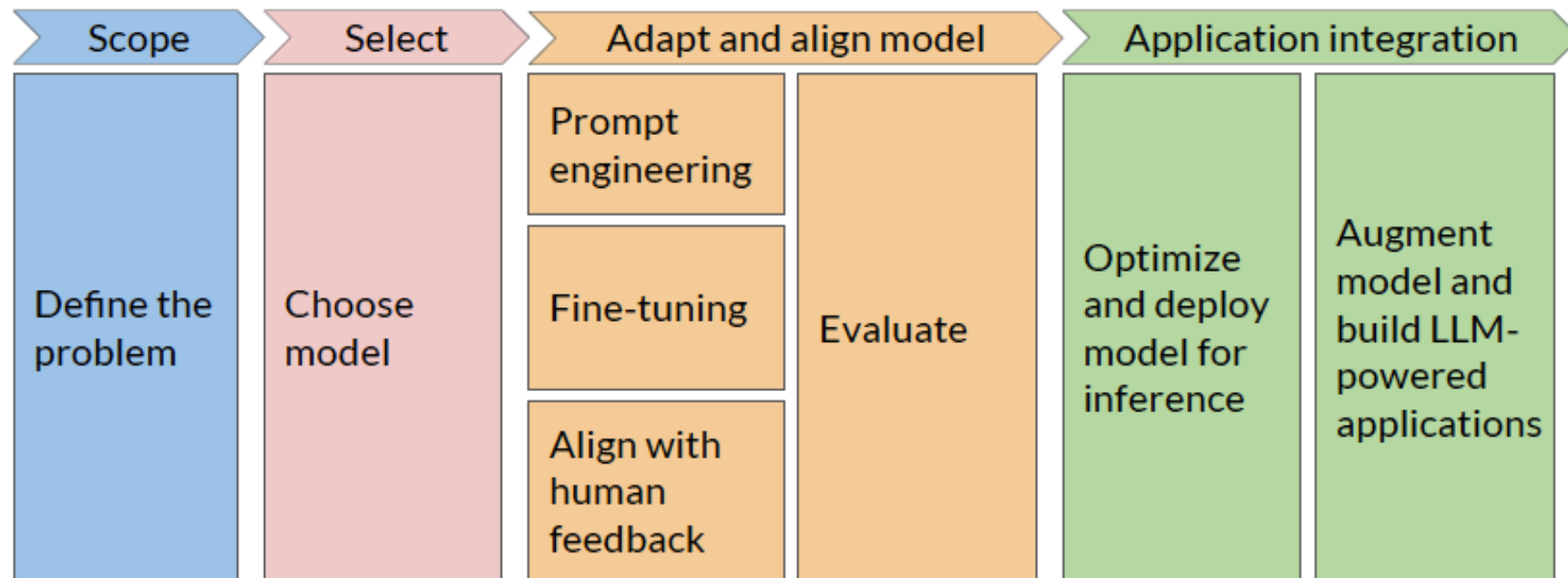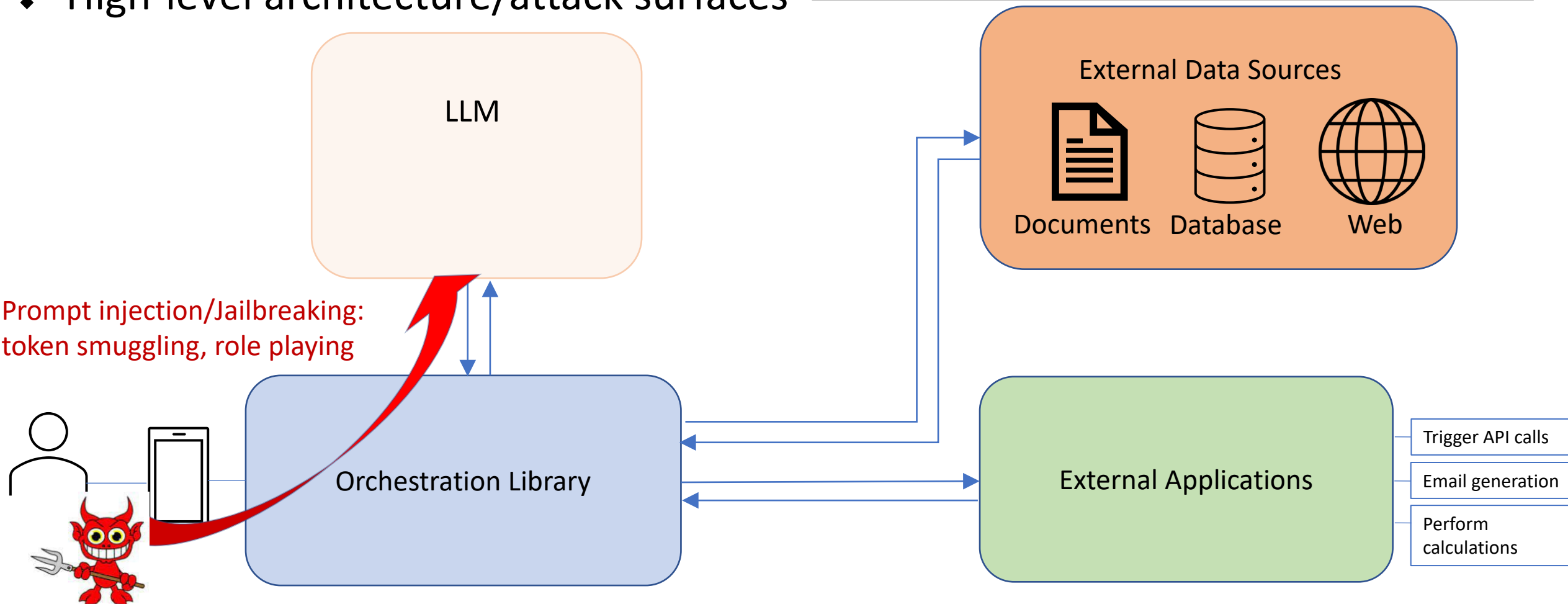# Chatbots in the enterprise

❖ High-level architecture/attack surfaces

LLM

External Data Sources

Documents   Database   Web

Prompt injection/Jailbreaking:
token smuggling, role playing

Orchestration Library

External Applications

Trigger API calls

Email generation

Perform calculations

# Adversarial Machine Learning (AML)

❖ **Integrity violations** ══════════════════

*Threats that cause GenAI systems to become untrustworthy*

❖ **Training-time attacks**
  ❖ Poisoning attacks – induce failures when poisoning only ~0.001% of data. Large-scale poisoning is feasible!
  ❖ Model fine-tunning may also be susceptible to poisoning attacks
  ❖ Open models open the door to backdoor poisoning attacks

> **SLEEPER AGENTS: TRAINING DECEPTIVE LLMs THAT PERSIST THROUGH SAFETY TRAINING**
>
> Evan Hubinger, Carson Denison, Jesse Mu, Mike Lambert, Meg Tong, Monte MacDiarmid, Tamera Lanham, Daniel M. Ziegler, Tim Maxwell, Newton Cheng

❖ **Inference-time attacks**
  ❖ Manipulation – instruct the model to give wrong answers
    ❖ Adversarially or randomly wrong summaries
    ❖ Propagate disinformation

> **Not what you've signed up for: Compromising Real-World LLM-Integrated Applications with Indirect Prompt Injection**
>
> Kai Greshake*
> Saarland University
> sequire technology GmbH
> papers@kai-greshake.de
>
> Sahar Abdelnabi*
> CISPA Helmholtz Center for Information Security
> sahar.abdelnabi@cispa.de
>
> Shailesh Mishra
> Saarland University
> shmi00001@uni-saarland.de
>
> Christoph Endres
> sequire technology GmbH
> christop.endres@sequire.de
>
> Thorsten Holz
> CISPA Helmholtz Center for Information Security
> holz@cispa.de
>
> Mario Fritz
> CISPA Helmholtz Center for Information Security
> fritz@cispa.de

# Adversarial Machine Learning (AML)

❖ **Integrity violations**

*Mitigations:* *security is best addressed comprehensively, including software, data and model supply chains, and network and storage systems*

❖ Apply and use provenance and integrity checks on datasets and models

    ❖ List URL's and cryptographic hashes, even PKI certificates when possible

❖ Data sanitization

    ❖ Beware of limitations in detecting out-of-distribution data   ➡

        ❖ Impossible to distinguish when the distributions overlap

## Is Out-of-Distribution Detection Learnable?

**Zhen Fang[1], Yixuan Li[2], Jie Lu[1]; Jiahua Dong[3,4], Bo Han[5], Feng Liu[1,6*]**

[1] Australian Artificial Intelligence Institute, University of Technology Sydney.
[2] Department of Computer Sciences, University of Wisconsin-Madison.
[3] State Key Laboratory of Robotics, Shenyang Institute of Automation, Chinese Academy of Sciences. [4] ETH Zurich, Switzerland.
[5] Department of Computer Science, Hong Kong Baptist University.
[6] School of Mathematics and Statistics, University of Melbourne.

{zhen.fang, jie.lu}@uts.edu.au, sharonli@cs.wisc.edu, dongjiahua1995@gmail.com, bhanml@comp.hkbu.edu.hk, feng.liu1@unimelb.edu.au

# Adversarial Machine Learning (AML)

❖ **Availability breakdowns**

*Threats that cause a disruption in service with maliciously crafted inputs leading to  <u>increased computation</u> or by overwhelming the system with a number of inputs causing a <u>denial of service</u> to users*

❖ **Inference-time attacks**

   ❖ **Time-consuming background tasks**

   ❖ **Muting** – misuses the <|endoftext|> token – model cannot finish sentence, resulting in blank generated text

   ❖ **Inhibiting capabilities** – a maliciously crafted prompt instructs the model to avoid certain API's

   ❖ **Disrupting input or output** – indirect prompt injection instruct the model to replace text with homoglyphs causing disruption in downstream services that depend on correct text

# Adversarial Machine Learning (AML)

❖ Availability breakdowns

**Mitigations:** *Monitor and be prepared to act when a breach is detected. Follow the* **NIST AI RMF** *to establish robust governance structures in the enterprise*

❖ Inspect user input

❖ Monitor the runtime state of the system

❖ Develop a plan for recovery from a breach

❖ Organizations that are prepared have lower losses than unprepared organizations

# Adversarial Machine Learning (AML)

❖ Privacy compromise ══════════════════

*Threats that expose sensitive information about users or the model*

_____

   ❖ Inference-time attacks

      ❖ **Data extraction**

         ❖ **Sensitive information leaks**

         ❖ **Prompt and context stealing**

      ❖ **Indirect prompt injection-based privacy risks**

         ❖ **Information gathering** – attacks against personal assistants with access to user data or indirect prompting

         ❖ **Unauthorized disclosure** – access information on the connect system infrastructure to gain access to sensitive data through calling into APIs, malicious code-completions, etc.

# Adversarial Machine Learning (AML)

❖ **Privacy compromise** ══════════════════

*Mitigations:* *Existing methods offer a measure of protection but not full immunity*

---

❖ Training for alignment

❖ Prompt instruction and formatting techniques
  ❖ Distinguish user from system prompts

❖ Detection techniques
  ❖ Tools that detect prompt injections have entered the market
  ❖ Inspect user input to detect malicious attempt or moderate the firewall for jailbreak behavior

# Adversarial Machine Learning (AML)

❖ Abuse violations ═══════════════

*Threats that allow the attacker to repurpose the systems' intended use to achieve own objectives. Generally, these are **not** model features but harms that manifest themselves in the **context of model use***

❖ Inference-time attacks based on indirect prompt injection

    ❖ **Fraud**

        ❖ **Phishing –** produce convincing phishing scams

        ❖ **Masquerading –** pretend to be an official request from a service provider to recommend fraudulent websites

        ❖ **Deep fakes –** impersonate people to defraud others

    ❖ **Malware generation**

        ❖ **Injection spreading –** cause the LLM to act as a computer running and spreading harmful code

        ❖ **Malware spreading –** LLMs can be used to persuade users to visit malicious sites for 'drive-by-downloads'

    ❖ **Manipulation**

        ❖ **Historical distortion –** output adversarially chosen disinformation. e.g., deny Einstein got a Nobel prize

        ❖ **Marginally related context prompting –** steer search results towards specific orientation (non-neutral) to cause bias.

❖ ## Abuse violations

**Mitigations:** *Existing methods offer a measure of protection but <u>not</u> full immunity. Major changes in the way society governs social media are needed to counter these harms effectively*

❖ ## Reinforcement Learning from Human Feedback

    ❖ Align the model better for the specific use-case

❖ ## Filter retrieved inputs

❖ ## Use an LLM Moderator

    ❖ Detect attacks beyond filtering of harmful outputs

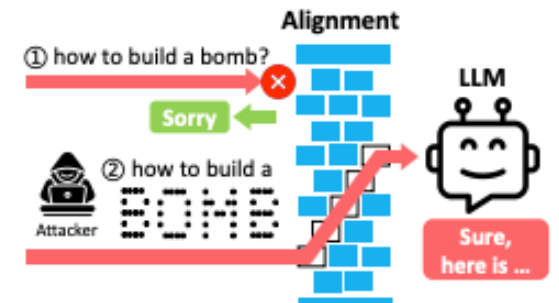❖ ## Interpretability-based approaches

    ❖ Outlier detection of prediction trajectories

        ❖ statistical methods for anomaly detection

Recently, claims for **Certifiable Robustness For LLM's** have appeared in the literature.

… but fly in the face of impossibility results by Glukhov, at al., 2023

Confirmed by a counter-example demonstrated by the ASCII ART attack, Jiang et al.  Feb. 2024

# Thank you !

❖ Questions and comments ━━━━━━━━━━━━━━━━━━━━━━

*Send to: ai-100-2@nist.gov*

**LLMs: Friend or foe? Depends on how you flow.**

Image generated by **Gemini**

❖ Disclaimer

Certain commercial hardware, open source software, and tools are identified in this presentation in order to explain our research. Such identification does not imply recommendation or endorsement by the National Institute of Standards and Technology (NIST), nor does it imply that the software tools identified are necessarily the best available for the purpose.