



**THALES**  
Building a future we can all trust

# Safeguarding ML: A comprehensive security plan

**Viswanath S CHIRRAVURI**  
October 2024



[www.thalesgroup.com](http://www.thalesgroup.com)



# Artificial Intelligence... What are we talking about?

## ➤ Artificial Intelligence

Any technique that enables computers to mimic human intelligence

## ➤ Machine Learning

A subset of Artificial Intelligence focusing on data-based learning by opposition to symbolic AI (aka knowledge-based learning).

## ➤ Deep Learning

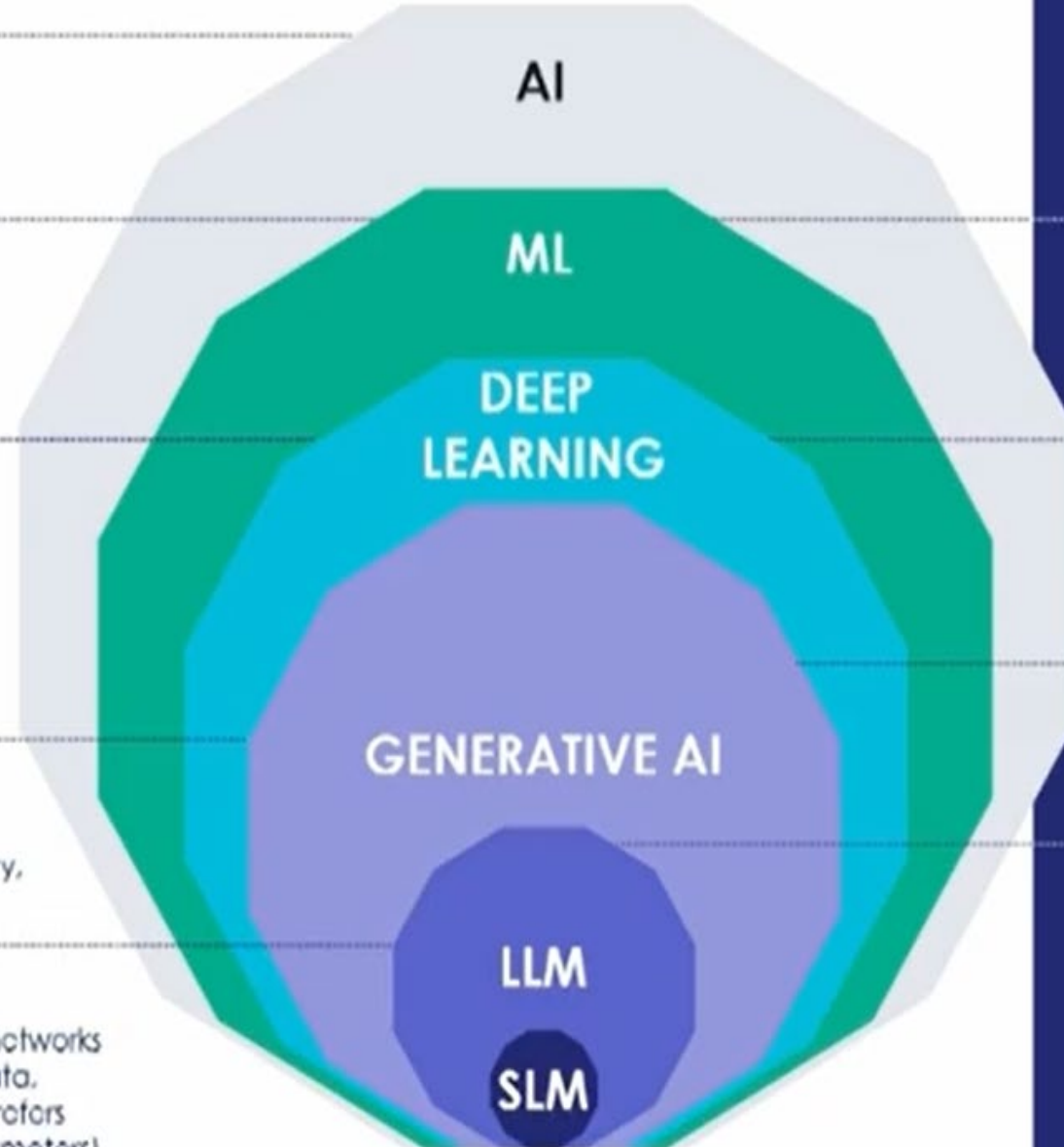
A subset of Machine Learning methods, based on Artificial Neural Networks with many layers (aka Deep Neural Networks). There are many types of ANN (e.g. CNN: Convolutional Neural Network, RNN: Recurrent Neural Network).

## ➤ Generative AI

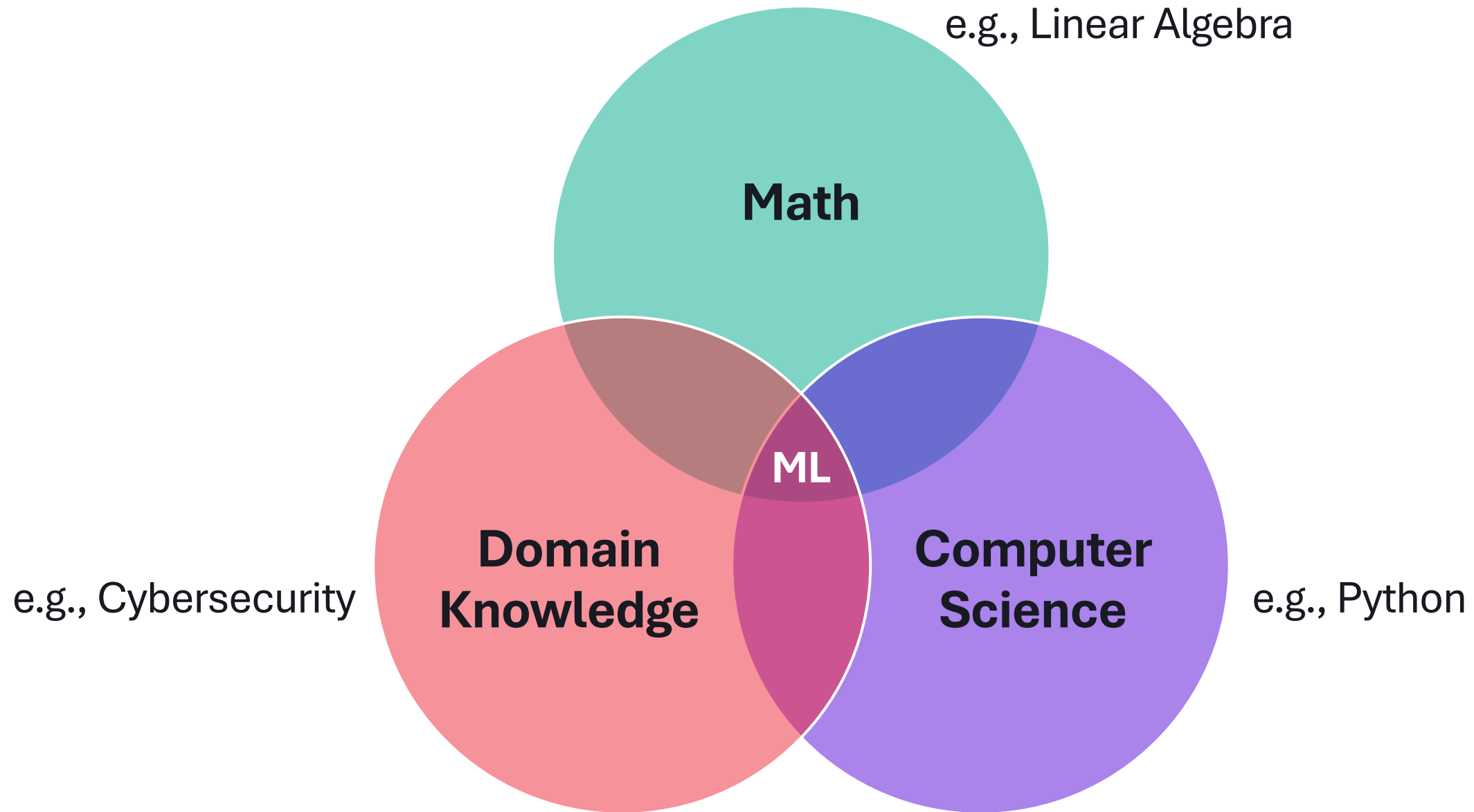
A type of artificial intelligence technology that can produce various types of content, including text, imagery, audio and synthetic data. Examples: GAN\*, LLM...

## ➤ Large Language Models

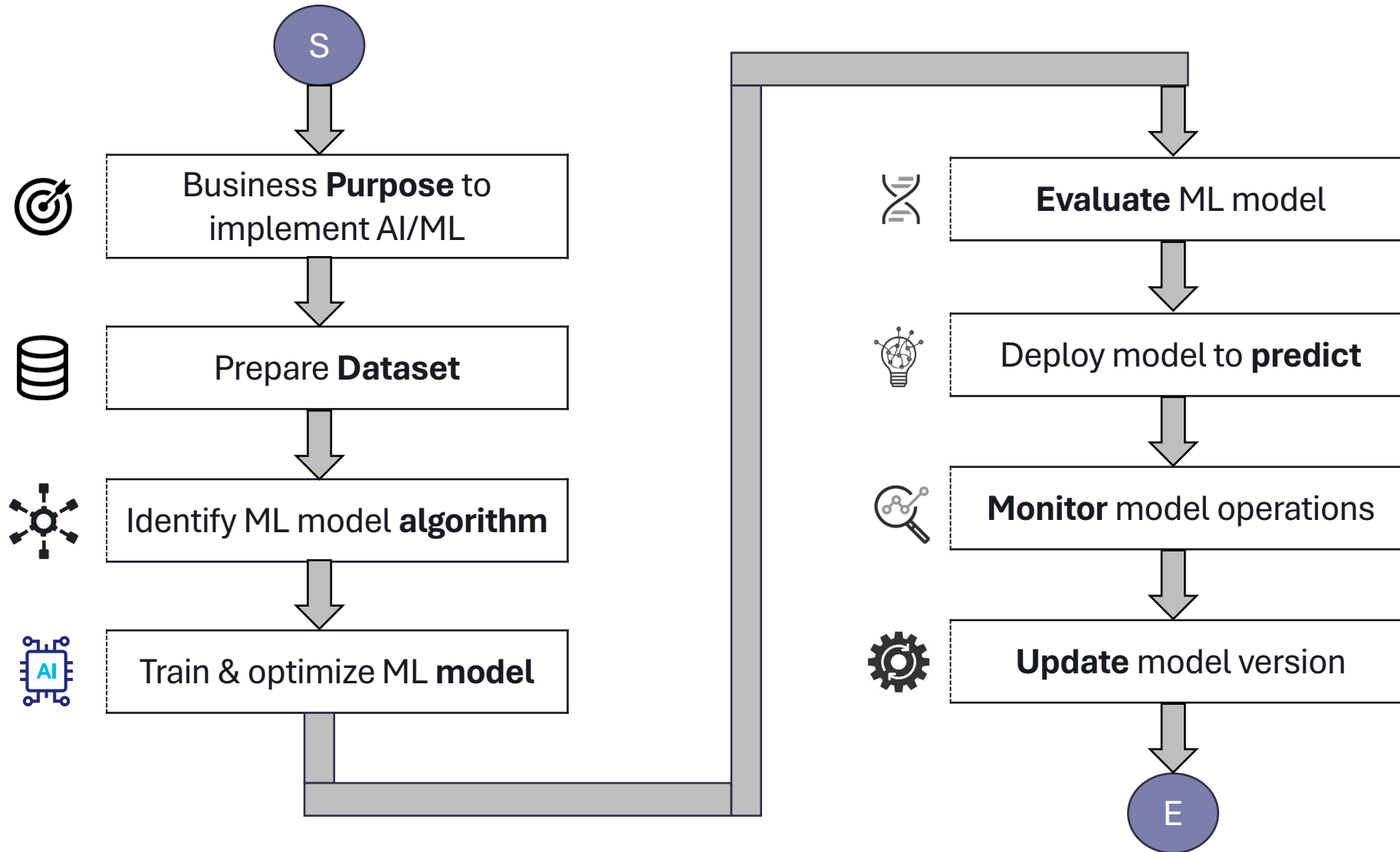
Specific application of GenAI using a variety of neural networks (the transformers), trained on a very large amount of data, commonly coming with XX+ billions of parameters. SLM refers to Small Language Model (in millions or X billions of parameters).



\*GAN: Generative Adversarial Network

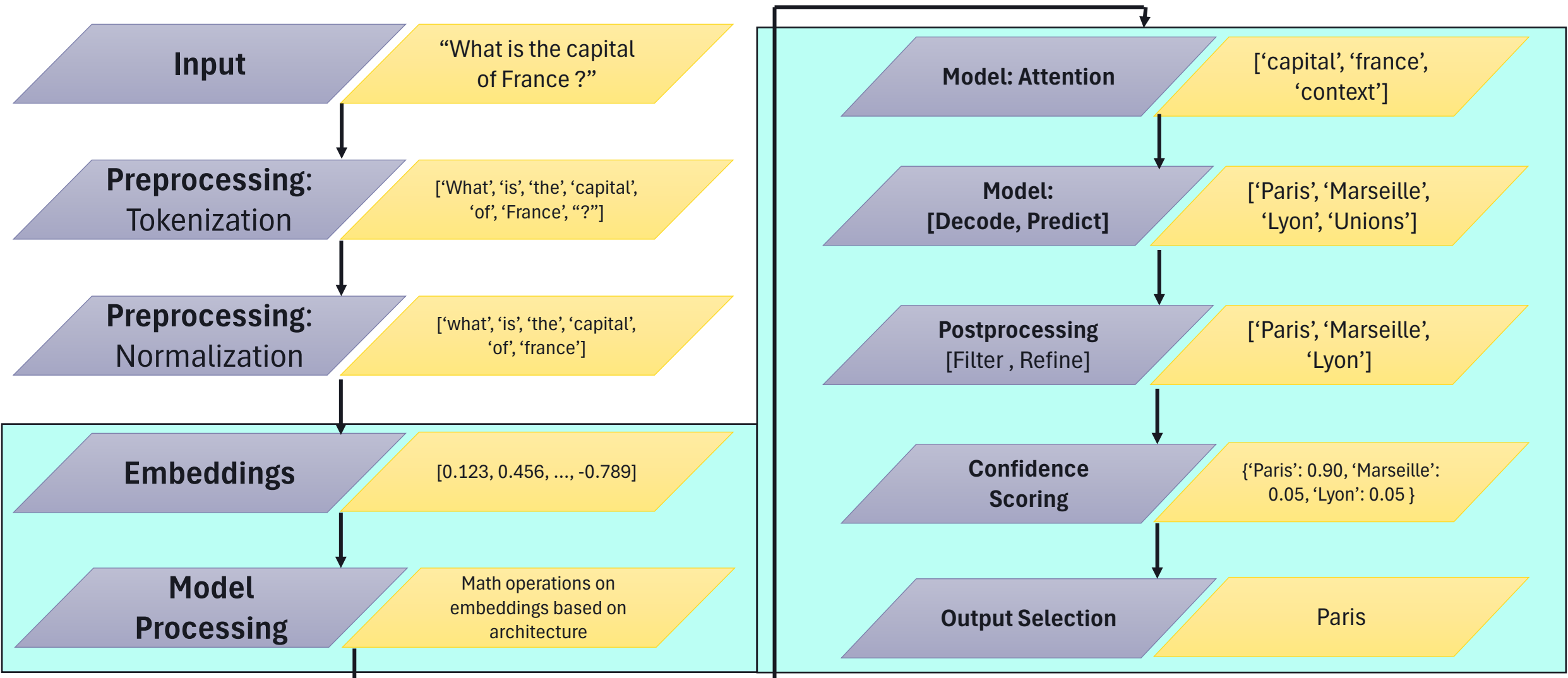


# Machine Learning Workflow



OPEN

# What happens when you input something into LLMs in production? (with example)





# Artificial Intelligence and Cybersecurity

## AI for Cybersecurity

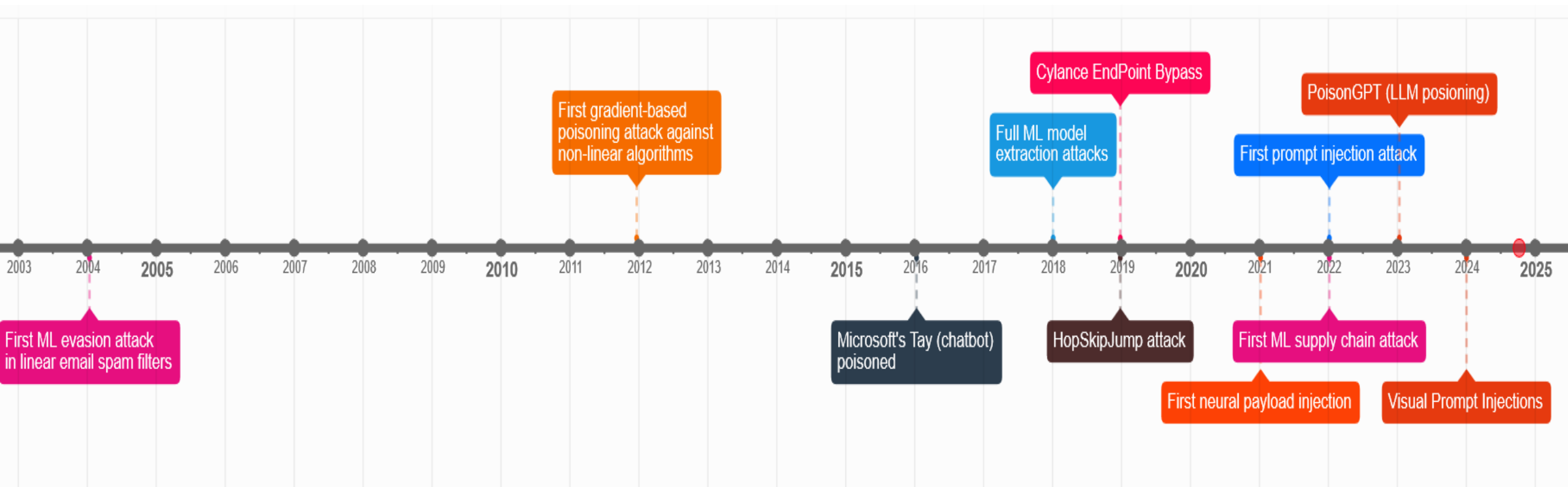
- Implementing AI/ML to protect systems (defense) e.g., AI for improved Data Protection
- Implementing AI/ML to hack systems (attack) e.g., AI generated malware, AI to scan for vulns
- Implementing AI/ML in forensics or incident handling (resolve) e.g., AI for threat hunting
- Implementing AI/ML in management (report) e.g., AI for Cybersecurity metrics / KPIs

## Cybersecurity for AI

- Implementing security in AI/ML business usecases (Trusted AI/ML) e.g., Securing LLMs

# Evolution of Cyberattacks in AI/ML

Source: HiddenLayer



<https://oecd.ai/en/incidents>  
<https://airisk.io/>  
<https://avidml.org/database/>  
<https://incidentdatabase.ai/>  
<https://atlas.mitre.org/studies>

OPEN

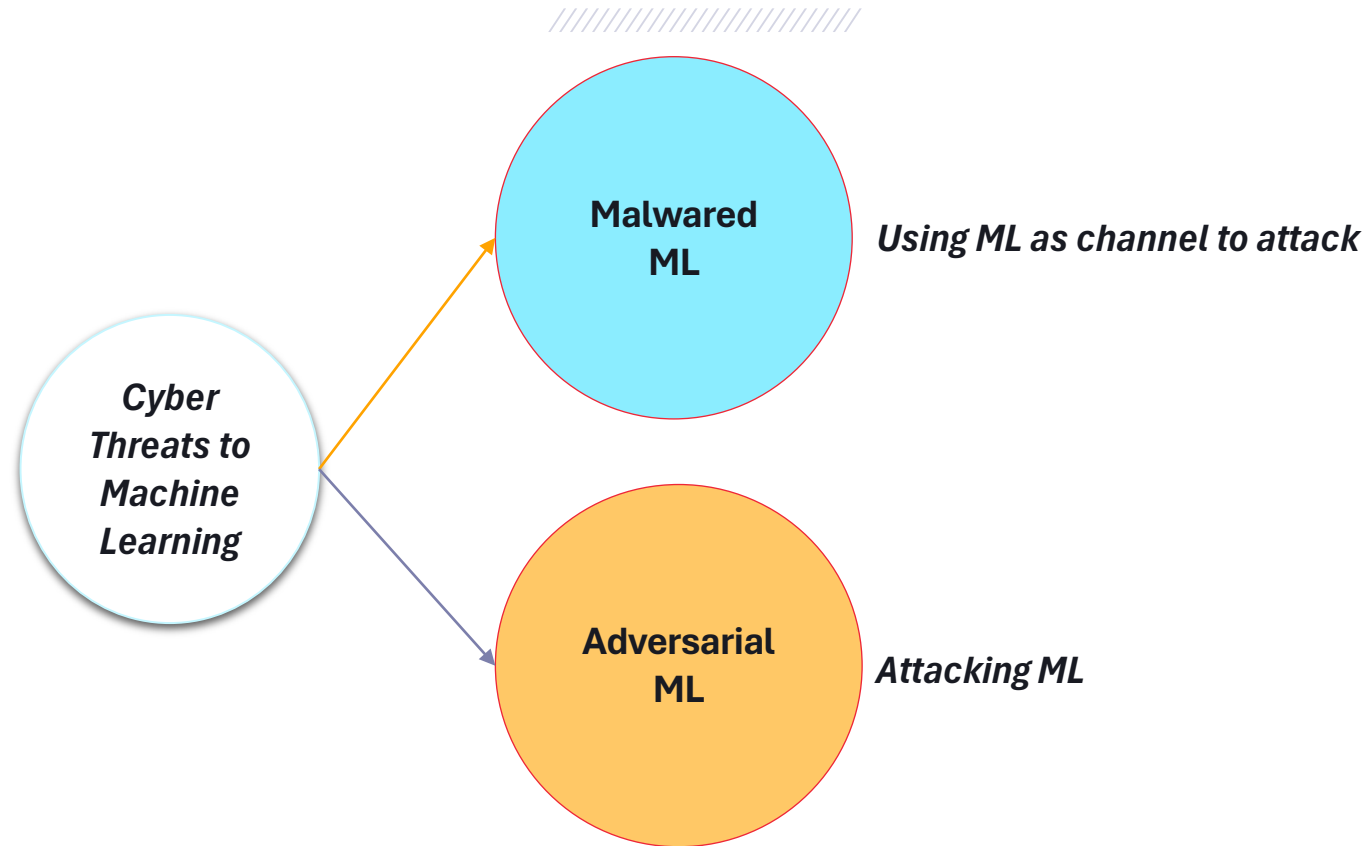
## BUSINESS RISKS TO IMPLEMENT AI / ML

- > Harmful Content Creation
- > Deepfakes
- > Copyright Violation
- > Reputational Damage
- > Increased Costs
- > Regulatory Fines
- > Intellectual Property Loss
- > Business Continuity Disruption
- > Customer Loss





# CYBER THREAT CATEGORIES TO MACHINE LEARNING



OPEN

# Cyber Threats to Machine Learning

## Malware ML

ML datasets, pre-trained models, or underlying ML libraries like TensorFlow, PyTorch, scikit-learn, etc. containing **malware** to spread across networks.

- ▶ Virus
- ▶ Worms
- ▶ Trojan horse
- ▶ Spyware
- ▶ Adware
- ▶ Botnets
- ▶ Ransomware

## Adversarial ML

Manipulate ML datasets, models, or systems with the goal of compromising the confidentiality, integrity, or availability of the business use case.

- ▶ Model Evasion
- ▶ Model Extraction / Theft
- ▶ Model Backdoor
- ▶ Model Denial of service
- ▶ Data Poisoning / Data Backdoor
- ▶ Data Theft / Model Inference
- ▶ Generative AI: Prompt / Code Injection / Jailbreaking

# AI/ML: Regulations, Standards, Frameworks

## Regulations

- US White House EO 14110
- EU AI Act
- Canada AI & Data Act
- UK AI Bill
- PDPC Singapore

## Standards

- NIST AI RMF
- ENISA & ETSI
- OWASP Top 10 (ML & LLM)
- ISO 27090
- IEEE P3119

## Frameworks

- MITRE ATLAS
- Gartner AI TRiSM
- Linux Foundation AIOSS
- Google SAIF
- IBM Secure Gen AI
- Thales Secure ML

# ThalesGroup's **Secure Machine Learning** Framework

> We open sourced this secure-machine-learning framework

▸ <https://github.com/ThalesGroup/secure-ml> [License: CC BY-ND 4.0]

> Includes:

- ❖ Corporate Security Policy: Framework for ML systems [[Link](#)]
- ❖ Security Requirements & Guidelines [[Link](#)]
- ❖ Privacy-Preserving Techniques [[Link](#)]
- ❖ Recommended Security Tools [[Link](#)]
- ❖ Industry references on regulations, standards [[Link](#)]
- ❖ Cyber threats to ML & Taxonomy [[Link](#)]

# ML Security Policy – highlights



- **Data Security:** It refers to securing data that can be training data, validation data, data stored in vector databases, etc.
- **Model Security:** It refers to hardening the model in both development and inference time.
- **Platform Security:** It refers to hardening of the underlying platform where the ML model is being developed or operated, along with any associated data.
- **Security Compliance:** It refers to being compliant to both internal and external regulations, standards, and best practices as application to the Trusted AI principles.
- **Human Security:** It refers to humans being well trained, authenticated, authorized, and capable to handle security incidents around ML ecosystem.



	DATA SECURITY	MODEL SECURITY	PLATFORM SECURITY	SECURITY COMPLIANCE	HUMAN SECURITY
GOAL	<ul style="list-style-type: none"><li>Confidentiality</li><li>Integrity</li><li>Availability</li><li>Authentication</li><li>Authorization</li><li>Non-repudiation</li><li>Privacy</li></ul>	<ul style="list-style-type: none"><li>Integrity in computation</li><li>Accuracy and precision in output</li></ul>	<ul style="list-style-type: none"><li>API security</li><li>System security</li><li>Network security</li></ul>	<ul style="list-style-type: none"><li>Comply with internal and external policies, standards and regulations</li></ul>	<ul style="list-style-type: none"><li>People involved are aware of security risks</li><li>Only authorized people are involved</li></ul>
TECHNIQUE	<ul style="list-style-type: none"><li>Encryption</li><li>Access Controls</li><li>Data backups &amp; Recovery</li><li>Data anonymization</li><li>Data quality control</li><li>Secure data sharing</li><li>Data classification</li></ul>	<ul style="list-style-type: none"><li>Secure development</li><li>Input and Output validation</li><li>Model explainability</li><li>Adversarial training</li><li>Robustness testing</li><li>Monitoring and alerting</li><li>Secure deployment</li></ul>	<ul style="list-style-type: none"><li>Vulnerability scanning &amp; Penetration Testing</li><li>Patch management</li><li>Access Controls</li><li>Encryption</li><li>Hardening</li><li>Secure Configuration</li></ul>	<ul style="list-style-type: none"><li>Regulatory compliance</li><li>Ethical considerations</li><li>Data retention and deletion</li><li>Audit trails</li><li>Security assessments</li><li>Third-party risk management</li></ul>	<ul style="list-style-type: none"><li>Training and awareness</li><li>Background checks</li><li>Incident response</li><li>Governance and oversight</li><li>Continuous monitoring</li></ul>



# ML Security Requirements – DATA SECURITY

WHAT?

## Data Security Techniques

ID	Technique	Activity
DS1	Encryption	<ul style="list-style-type: none"><li>- Encrypt sensitive data in transit and at rest.</li><li>- Use strong algorithms (e.g., AES, RSA).</li><li>- Securely manage encryption keys.</li></ul>
DS2	Access Controls	<ul style="list-style-type: none"><li>- Implement strict access controls.</li><li>- Ensure authorized personnel access only.</li><li>- Utilize user authentication, role-based access, and MFA.</li></ul>
DS3	Data Classification, Backups, and Recovery	<ul style="list-style-type: none"><li>- Classify data by sensitivity.</li><li>- Regularly back up data.</li><li>- Develop a disaster recovery plan.</li></ul>
DS4	Data Anonymization	<ul style="list-style-type: none"><li>- Remove personally identifiable information (PII).</li><li>- Use data masking or tokenization techniques to protect privacy.</li></ul>
DS5	Data Quality Control	<ul style="list-style-type: none"><li>- Implement processes for data quality and integrity.</li><li>- Conduct data validation and antivirus scans.</li></ul>
DS6	Secure Data Sharing	<ul style="list-style-type: none"><li>- Ensure security measures for third-party data sharing.</li><li>- Use secure transfer protocols and establish data usage agreements.</li></ul>

## Key Takeaways

1. Implement Robust Security Measures: Focus on encryption, access controls, and data quality.
2. Protect Sensitive Information: Anonymization and secure sharing are essential.
3. Plan for Recovery: Regular backups and disaster recovery plans ensure data integrity.

# ML Security Requirements – MODEL SECURITY

WHAT?

## Model Security Techniques

ID	Technique	Activity
MS1	Secure Development	<ul style="list-style-type: none"><li>- Implement version control for AI models.</li><li>- Use code repositories and branching.</li><li>- Conduct code reviews and maintain a secure development environment.</li></ul>
MS2	Input/Output Validation	<ul style="list-style-type: none"><li>- Validate inputs and outputs for accuracy and reliability.</li><li>- Employ data validation, normalization, and transformation to eliminate untrusted inputs.</li></ul>
MS3	Model Explainability	<ul style="list-style-type: none"><li>- Ensure AI models are explainable and transparent.</li><li>- Utilize feature importance analysis and model visualization.</li><li>- Conduct threat modeling.</li></ul>
MS4	Adversarial Training & Robustness Testing	<ul style="list-style-type: none"><li>- Perform robustness testing against malware and adversarial attacks.</li><li>- Implement scanning, adversarial training, and model hardening techniques.</li></ul>
MS5	Monitoring and Alerting	<ul style="list-style-type: none"><li>- Monitor AI models in production for anomalies.</li><li>- Set up logging and alerting mechanisms for unexpected behavior.</li></ul>
MS6	Secure Deployment	<ul style="list-style-type: none"><li>- Ensure secure deployment of AI models.</li><li>- Use secure communication protocols and access controls.</li><li>- Implement authentication and reverse engineering protection.</li></ul>

## Key Takeaways

1. Integrity and Resilience: Focus on secure development and robustness against attacks.
2. Transparency: Prioritize model explainability for better understanding and trust.
3. Continuous Monitoring: Regularly monitor models to detect anomalies and ensure secure deployment.

# ML Security Requirements – PLATFORM SECURITY

WHAT?

## Platform Security Techniques

ID	Technique	Activity
PS1	Vulnerability Scanning	<ul style="list-style-type: none"><li>- Perform vulnerability scanning and penetration testing.</li><li>- Identify and mitigate vulnerabilities in AI/ML infrastructure and systems.</li></ul>
PS2	Patch Management	<ul style="list-style-type: none"><li>- Implement a robust patch management process.</li><li>- Ensure timely installation of software and firmware updates to address known vulnerabilities.</li></ul>
PS3	Access Controls	<ul style="list-style-type: none"><li>- Enforce strict access controls for AI/ML infrastructure.</li><li>- Use user authentication, role-based access control, and multi-factor authentication.</li></ul>
PS4	Encryption	<ul style="list-style-type: none"><li>- Utilize encryption for data in transit and at rest.</li><li>- Implement TLS for secure communications and full-disk encryption for storage.</li></ul>
PS5	Network Hardening	<ul style="list-style-type: none"><li>- Apply network security measures such as firewalls and intrusion detection/prevention systems.</li><li>- Implement network segmentation to defend against attacks.</li></ul>
PS6	Hardware Security	<ul style="list-style-type: none"><li>- Secure hardware components (e.g., GPUs, TPUs) against tampering and physical attacks.</li></ul>
PS7	Secure Configuration	<ul style="list-style-type: none"><li>- Ensure secure configuration of ML systems.</li><li>- Follow best practices for hardening networks, systems, and applications, both in cloud and on-premises.</li></ul>

## Key Takeaways

1. Proactive Security: Regular vulnerability scanning and patch management are essential.
2. Controlled Access: Implement robust access controls to safeguard infrastructure.
3. Comprehensive Protection: Ensure both hardware and software configurations are secure.

# ML Security Requirements – SECURITY COMPLIANCE

WHAT?

## Security Compliance Techniques

ID	Technique	Activity
SC1	Regulatory Compliance	<ul style="list-style-type: none"><li>- Ensure AI/ML systems comply with regulations (e.g., HIPAA, GDPR, PCI-DSS, ISO 27001).</li><li>- Implement data protection, security controls, and privacy requirements.</li></ul>
SC2	Ethical Considerations	<ul style="list-style-type: none"><li>- Consider ethical implications, focusing on bias, fairness, and accountability.</li><li>- Employ bias mitigation techniques and promote algorithmic transparency.</li></ul>
SC3	Data Retention and Deletion	<ul style="list-style-type: none"><li>- Establish policies for data retention and deletion to meet regulatory requirements.</li><li>- Implement data minimization, anonymization, and destruction practices.</li></ul>
SC4	Audit Trails	<ul style="list-style-type: none"><li>- Create audit trails to track data access and usage.</li><li>- Utilize logging, monitoring, and reporting, especially in regulated industries.</li></ul>
SC5	Security Assessments	<ul style="list-style-type: none"><li>- Conduct security assessments to identify vulnerabilities and ensure compliance.</li><li>- Perform security risk assessments, audits, and compliance evaluations.</li></ul>
SC6	Third-Party Risk Management	<ul style="list-style-type: none"><li>- Ensure third-party vendors comply with regulatory requirements and security standards.</li><li>- Assess risks associated with vendors accessing sensitive data or systems.</li></ul>

## Key Takeaways

1. Compliance Matters: Adhere to regulations to protect data and privacy.
2. Ethical Frameworks: Address ethical considerations to ensure fairness and accountability.
3. Thorough Assessments: Regularly assess security and third-party risks for robust compliance.



## Human Security Techniques

ID	Technique	Activity
HS1	Training and Awareness	<ul style="list-style-type: none"><li>- Provide training programs for employees and stakeholders.</li><li>- Implement security awareness campaigns and conduct phishing simulations to highlight risks.</li></ul>
HS2	Background Checks	<ul style="list-style-type: none"><li>- Conduct background checks on employees and contractors with access to sensitive data.</li><li>- Focus on regulated industries and government agencies.</li></ul>
HS3	Incident Response	<ul style="list-style-type: none"><li>- Develop an incident response plan for handling security incidents.</li><li>- Include incident management, forensics, and communication plans with designated personnel.</li></ul>
HS4	Governance and Oversight	<ul style="list-style-type: none"><li>- Establish governance structures for secure and ethical AI/ML system development and deployment.</li><li>- Implement compliance frameworks and security dashboards.</li></ul>

### Key Takeaways

1. Empower Employees: Training and awareness are crucial for understanding security risks.
2. Vetting Access: Background checks enhance security for sensitive data access.
3. Preparedness: An effective incident response plan ensures readiness for security incidents.

HOW?

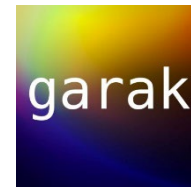
## ML Security Guidelines

	GUIDELINES
<b>DATA SECURITY</b>	OWASP ASVS, ISO27001, GDPR, IEEE, + custom guidelines
<b>MODEL SECURITY</b>	SHAP/LIME, ProtectAI' s LLM-guard, Microsoft AI Threat modeling + custom guidelines
<b>PLATFORM SECURITY</b>	Custom guidelines
<b>SECURITY COMPLIANCE</b>	FHE & FL (w/ multiple libraries), + custom guidelines
<b>HUMAN SECURITY</b>	Custom guidelines

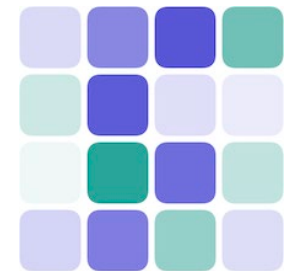


# Security tools for Machine Learning

Azure/**PyRIT**



QData/**TextAttack**



# Privacy-Preserving Techniques for Machine Learning

## Differential Privacy for Machine Learning

- <https://github.com/OpenMined/PyDP>
- <https://github.com/tensorflow/privacy>
- <https://github.com/opencv/opencv>

## Federated Learning

- <https://github.com/tensorflow/federated>
- <https://github.com/OpenMined/PySyft>
- <https://flower.dev/>

## Homomorphic Encryption

- <https://github.com/OpenMined/TenSEAL>
- <https://github.com/zama-ai/concrete-ml>
- SEAL (Simple Encrypted Arithmetic Library): <https://www.microsoft.com/en-us/research/project/microsoft-seal/>
- <https://github.com/homeinc/HElib>
- More at <https://github.com/jonasc hn/awesome-he>

## Secure Multi-Party Computation (SMPC)

- <https://github.com/data61/MP-SPDZ>
- <https://sharemind.cyber.ee/sharemind-mpc/>
- <https://github.com/easy-smpc/easy-smpc>

## Privacy-Preserving Data Synthesis

- <https://github.com/pytorch/opacus>
- <https://github.com/sdv-dev/SDV>
- <https://github.com/sdv-dev/CTGAN>
- <https://github.com/DataResponsibly/DataSynthesizer>

OPEN

# Secure ML Program Roadmap: 3, 6, 12-Month Plan

## 3

### Foundation

1. Draft/refine corporate security policy for ML systems.
2. Plan/red team AI/ML systems (create team if needed).
3. Analyze security regulations and standards for AI/ML
4. Raise awareness among ML engineers on security risks.

## 6

### Build

1. Research security tools for your ML tech stack.
2. Build guidelines for securing data pipelines and models.
3. Explore privacy-preserving techniques (e.g., FHE, FL).
4. Begin threat modeling for ML risks.

## 12

### Launch

1. Deliver secure ML solutions to clients.
2. Implement continuous monitoring for threats.
3. Automate security compliance for ML datasets & models.
4. Train stakeholders on new ML attack vectors.

# Learn by hands-on exercises

## 1. Exploring Adversarial Machine Learning by NVIDIA

- ▶ [https://learn.nvidia.com/courses/course-detail?course\\_id=course-v1:DLI+S-DS-03+V1](https://learn.nvidia.com/courses/course-detail?course_id=course-v1:DLI+S-DS-03+V1)

## 2. Adversarial Machine Learning

- ▶ <https://secml.readthedocs.io/en/v0.15/tutorials.adv.html>

## 3. Dreadnode's crucible (CTF)

- ▶ <https://crucible.dreadnode.io/>

## 4. Lakera's Gandalf (8-levels) for prompt injection attack (CTF)

- ▶ <https://gandalf.lakera.ai/baseline>



**Viswanath S CHIRRAVURI**

[https://www.linkedin.com/in/  
chviswanath/](https://www.linkedin.com/in/chviswanath/)



**THALES**  
Building a future we can all trust

# Thank you

[www.thalesgroup.com](http://www.thalesgroup.com)