



Overview **OWASP AI Exchange**

(owaspai.org)

Behnaz Karimi

13.11.2024



- Over 14 Years in Cybersecurity
- Senior Cybersecurity Analyst / AI Security Analyst at Accenture (AI-SDLC, Governance, Risk analysis)
- Passionate about AI Security/AI-Ransomware
- Co-author of GenAI Red Teaming Methodologies and best practices Top 10 LLM (still on process)
- Life time learner

OWASP and Me

- Been with AI Exchange for 1 year and 3 months
- Co-Lead /Lead AI Security Red Teaming/ Co-Author





German
OWASP
Day 2024



Rob van der Veer

Project leader of **OWASP AI Exchange**
/ **AI Privacy and Security Guide**



Chris Ancharski

Co-Lead



Aruneesh Salhotra

Co-Lead



Behnaz Karimi

Co-Lead



OWASP AI Exchange Mission



KEY FACTS

- Exchange is CC0 1.0 licensed: **free of copyright** and attribution.
- Volume is **170 pages**
- Accessible and Available at **owaspai.org**

The AI security community is marked with CC0 1.0 (Creative Common) meaning you can use any part freely, without attribution. If possible, it would be nice if the OWASP AI Exchange is credited and/or linked to, for readers to find more information.



Key Achievements

- Our direct flow of content into **ISO/IEC 27090** and standards for the **EU AI Act**.
 - **70 pages is accepted!**
- The liaison relation with CEN/CENELEC.
Established an official partnership between OWASP and CEN/CENELEC
Work for security starts from TC21/WG1 to TC21/WG5
- Recognition by standard makers across the globe
- Rob has been elected by the European countries as the co-editor of the official AI Act Security Standard.



Relation to other OWASP or other organization initiatives



The OWASP AI security and privacy guide is the official OWASP project under which the AI Exchange was established. The deliverables consists of the AI Exchange content plus guidance on AI privacy.

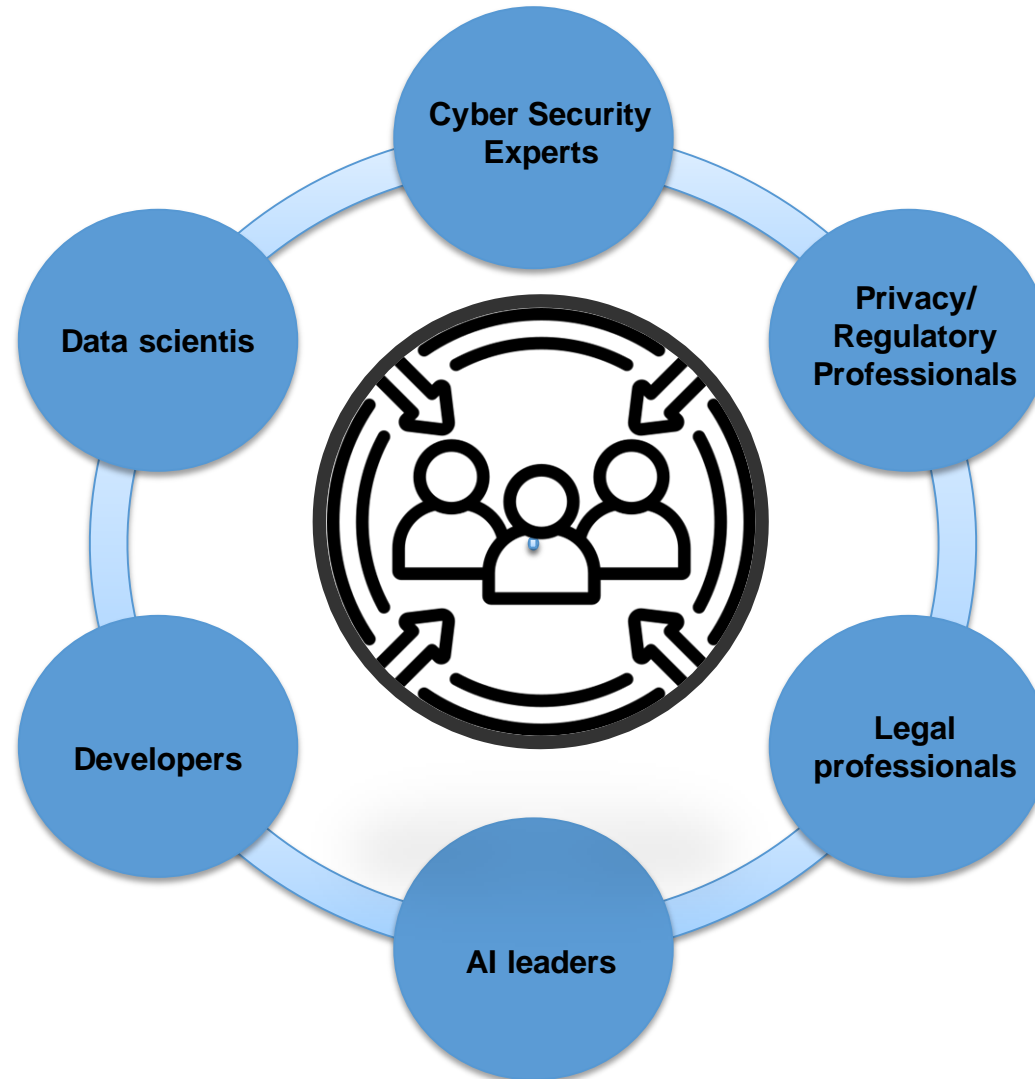
The OWASP LLM top 10 provides a list of the most important LLM security issues, plus deliverables that focus on LLM security, such as the LLM AI Security & Governance Checklist.

The OWASP ML top 10 provides a list of the most important machine learning security issues.

OpenCRE.org holds a catalog of common requirements across various security standards inside and outside of OWASP.



Target Audience





Scope & Responsibilities

- Develop a comprehensive framework for AI threats, risks, mitigations, and controls.
- Create a map integrating AI regulatory and privacy regulations.
- Establish a common taxonomy and glossary for AI security.
- Provide guidance on testing tools with outcome assessments.
- Formulate a shared responsibility model for third-party AI model usage.
- Offer supply chain guidance and an incident response plan.



Collaboration Efforts and Engagement

- We have regular collaboration with
 - CSA
 - ISO/IEC
- Liaison with CEN/CENELEC
- We have regular meetings with
 - NIST
 - MITRE
 - ITU
- We're part of the AISIC.





How to address AI Security?

Imperative to approach AI applications with a clear understanding of potential threats and which of those threats to prioritize for each use case.

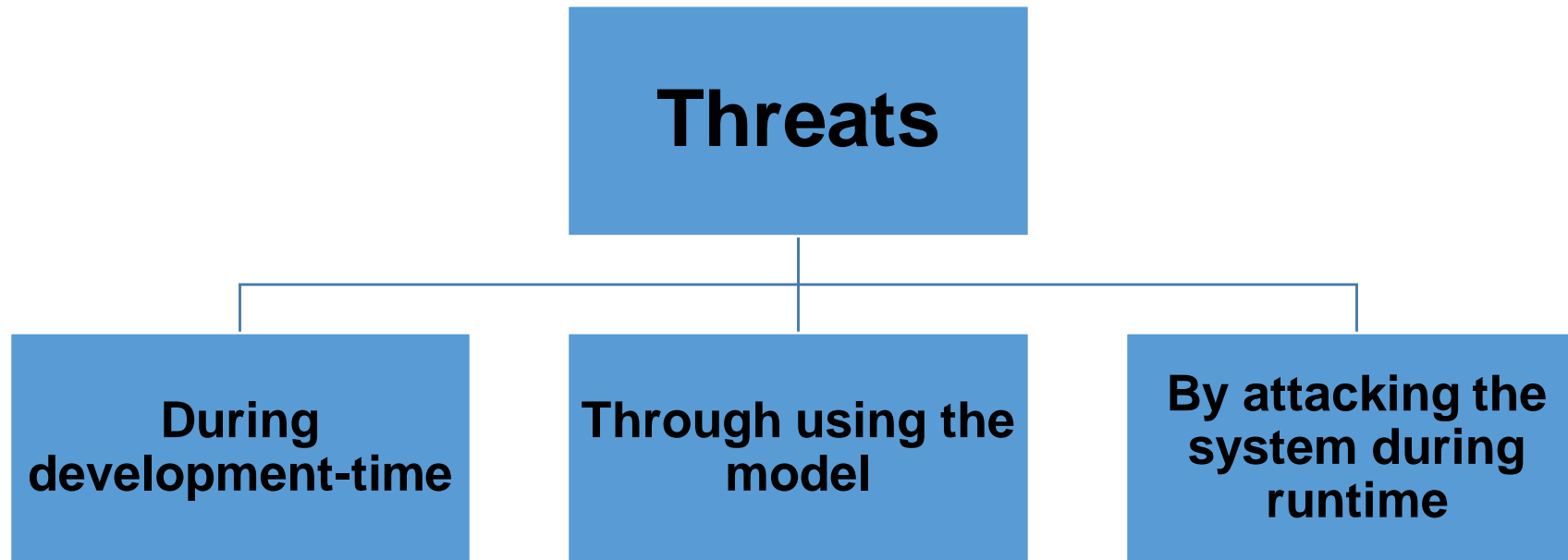
Standards and governance help guide this process for individual entities leveraging AI capabilities.

- *Implement AI governance*
- *Extend security and development practices*
- *Improve regular application and system security through understanding of AI particularities*

- *Limit the impact of AI by minimizing privileges and adding oversight*
- *Countermeasures in data science through understanding of model attacks*



Threat Model





Impacts

1. Confidentiality of **train/test data**
2. Confidentiality of **model Intellectual property**
(the *model parameters* or the process and data that led to them)
3. Confidentiality of **input data**
4. Integrity of **model behaviour** (the model is not manipulated to behave in an unwanted way)
5. **Availability** of the model
6. Confidentiality, integrity, and availability of **non AI-specific assets**



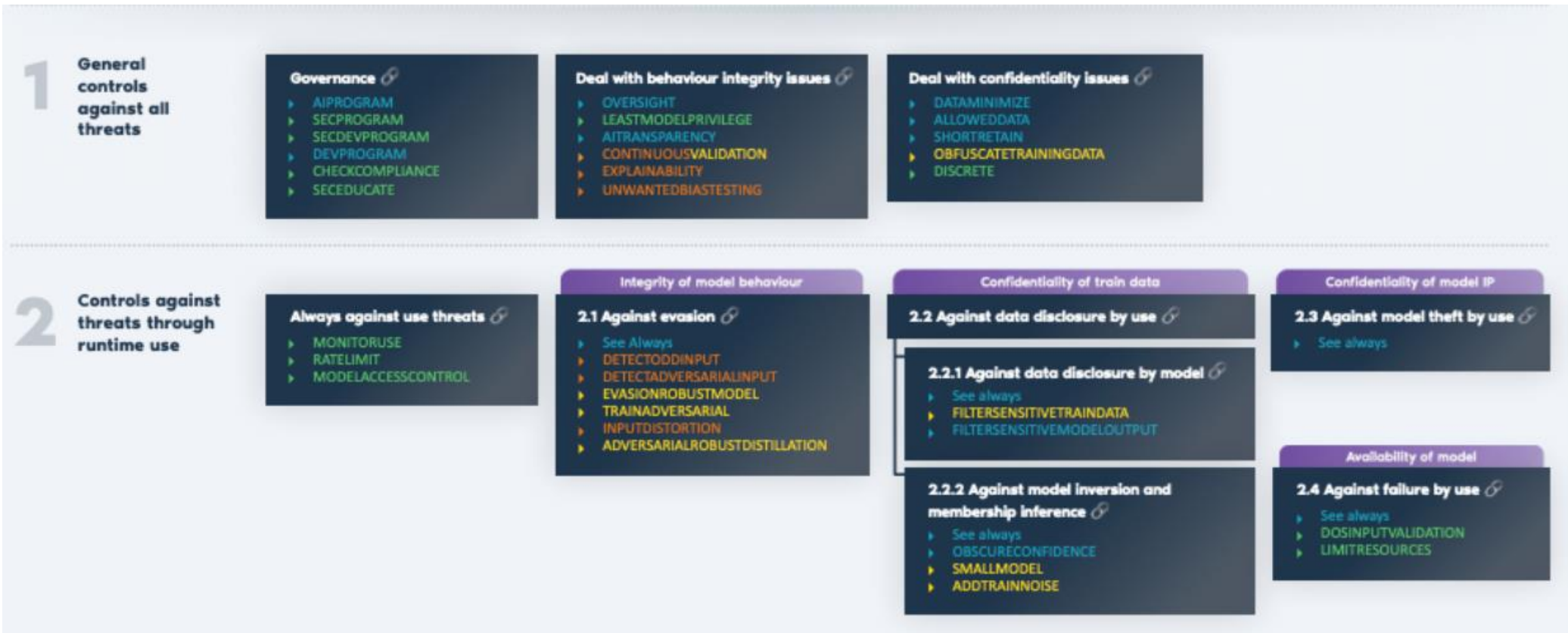
AI Security Matrix

AI-specific?	Lifecycle	Attack surface	Threat	Asset	Impacted	Unwanted result
AI	Runtime	Model use (provide input/ read output)	Direct prompt injection	Model behaviour	Integrity	Manipulated unwanted model behaviour causes wrong decisions leading to business financial loss, misbehaviour going undetected, reputational damage, legal and compliance issues, operational disruption, customer dissatisfaction and churn, reduced employee morale, incorrect strategic decisions, liability issues, personal damage and safety issues
			Indirect prompt injection			
			Evasion (e.g. adversarial examples)			
		Break into deployed model	Runtime model poisoning (reprogramming)			
	Development	Engineering environment	Development time model poisoning	Train data	Confidentiality	
		Supply chain	Data poisoning of train/finetune data			
			Obtain poisoned foundation model (transfer learning attack)			
			Obtain poisoned data to train/finetune			
	Runtime	Model use	Unwanted disclosure in model output	Train data	Confidentiality	
			Model inversion / Membership inference			
Development	Engineering environment	Train data leaks			Leaking sensitive data can cause costs from fines and legal fees and remediation effort, loss of business through customer churn, reputation damage, loss of competitive advantage in case of trade secrets, operational disruption, impacted business relationships, and employee morale	
Runtime	Model use	Model theft through by use (input-output harvesting)	Model intellectual property	Confidentiality	If attackers can copy a model, the investment in the model is devalued caused by loss of competitive advantage, plus a copy can help craft (evasion) attacks	
		Break into deployed model				Runtime model theft
Development	Engineering environment	Development time model parameter leak				
Runtime	Model use	System failure by use (model resource depletion)	Model behaviour	Availability	The model is not available, leading to business continuity issues, or safety problems	
Runtime	All IT		Model input leak	Model input data	Confidentiality	Sensitive data in model input leaks. E.g. an LLM prompt with a sensitive question, enhanced with retrieved company secrets
Runtime	All IT		Model output contains injection attack	Any asset	C, I, A	Injection attack (from model output) causes harm
Generic	Runtime	All IT	Generic runtime security attack	Any asset	C, I, A	Generic runtime security attack causes harm (includes social engineering/phishing)
	Development	All IT	Generic supply chain attack	Any asset	C, I, A	Generic supply chain security attack causes harm (e.g. vulnerability in a component)

Source: OWASP AI Exchange at owaspai.org



AI Security Threats and Controls Navigator



LEGEND:

Group of controls, ordered by threat or type (clickable)

Standard information security CONTROL (with attention points)

Runtime Data science CONTROL

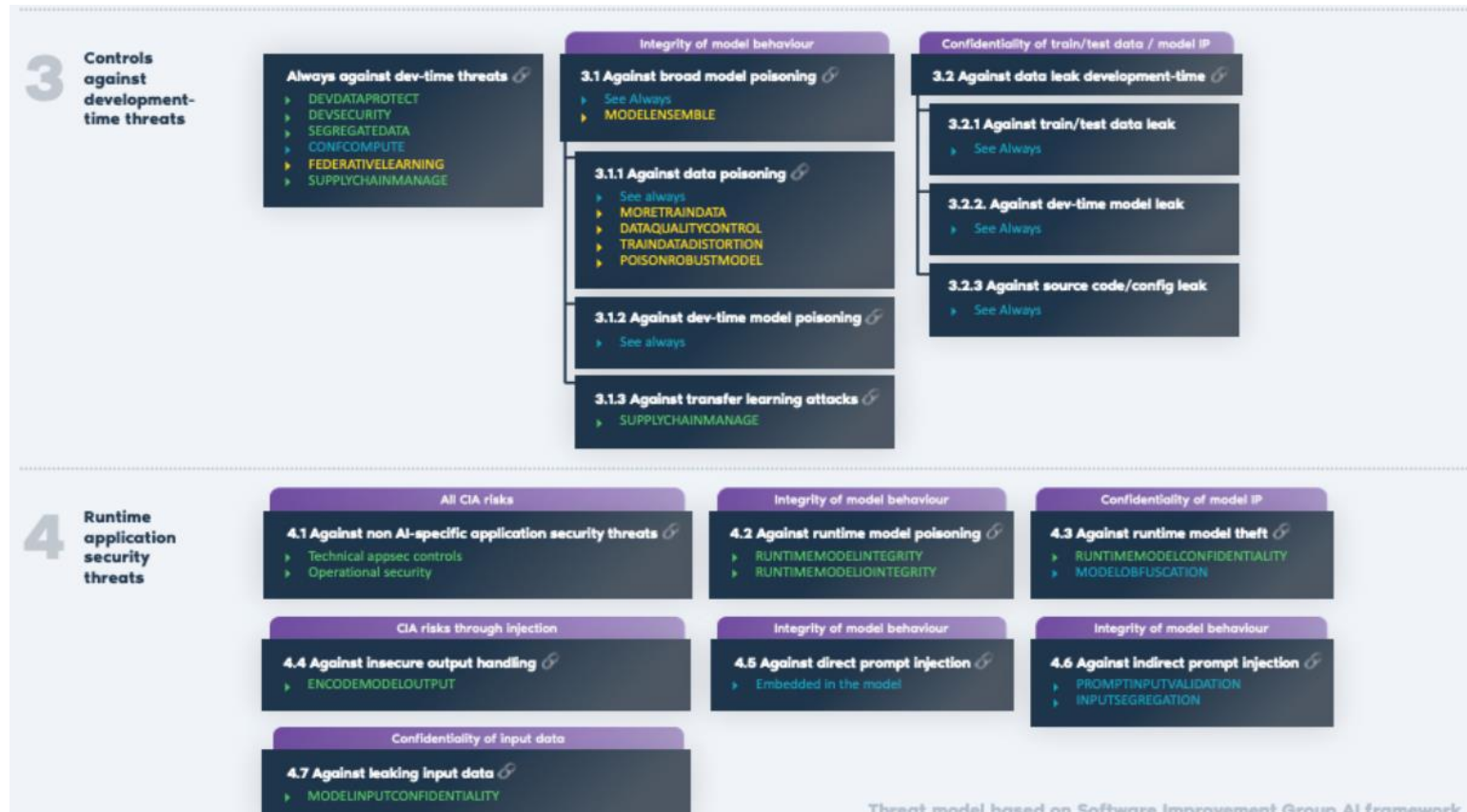
Development-time Data science CONTROL

Other CONTROL

Impact on Confidentiality, Integrity or Availability



AI Security Threats and Controls Navigator



LEGEND:

Group of controls, ordered by threat or type (clickable)

Standard information security CONTROL (with attention points)

Runtime Data science CONTROL

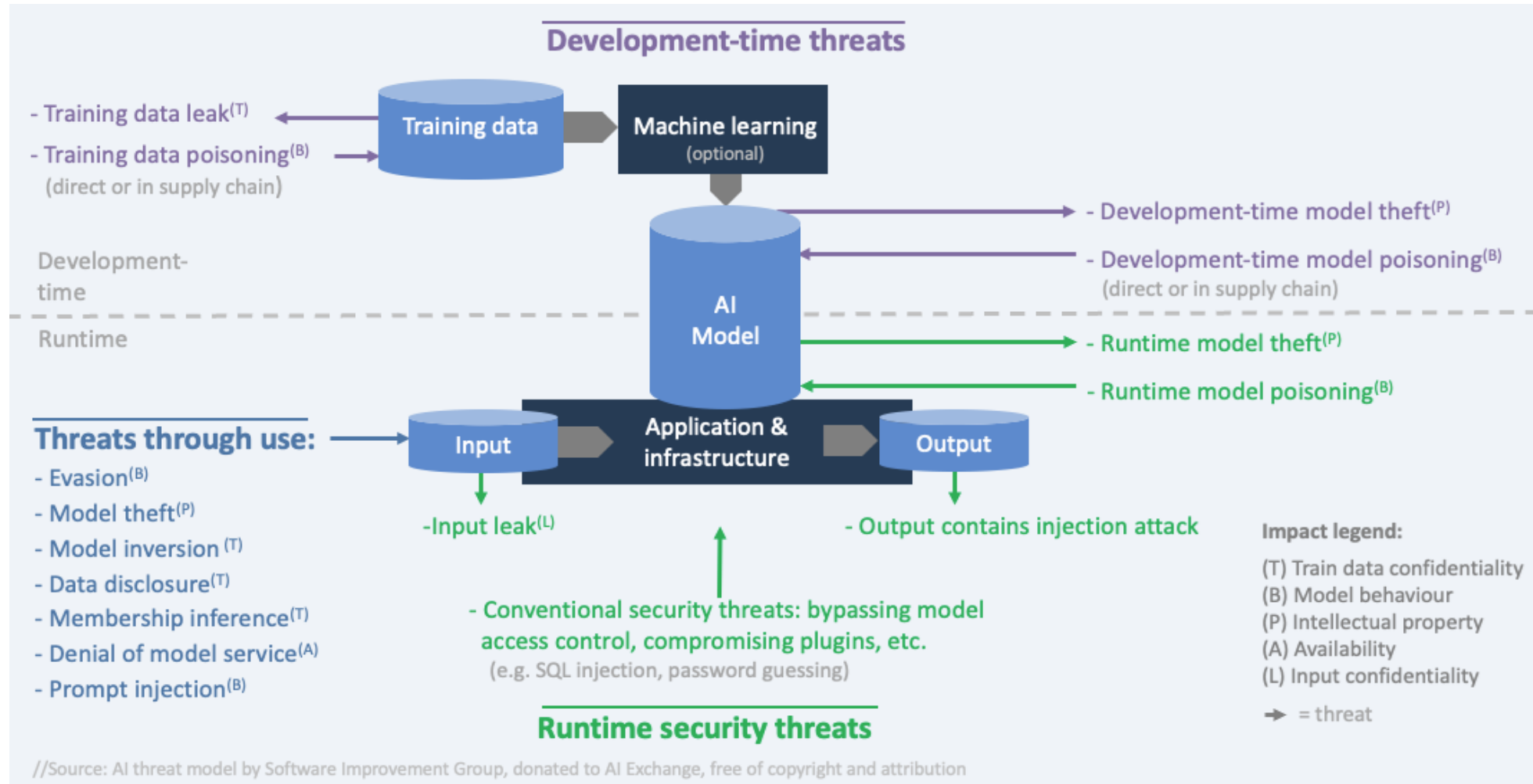
Development-time Data science CONTROL

Other CONTROL

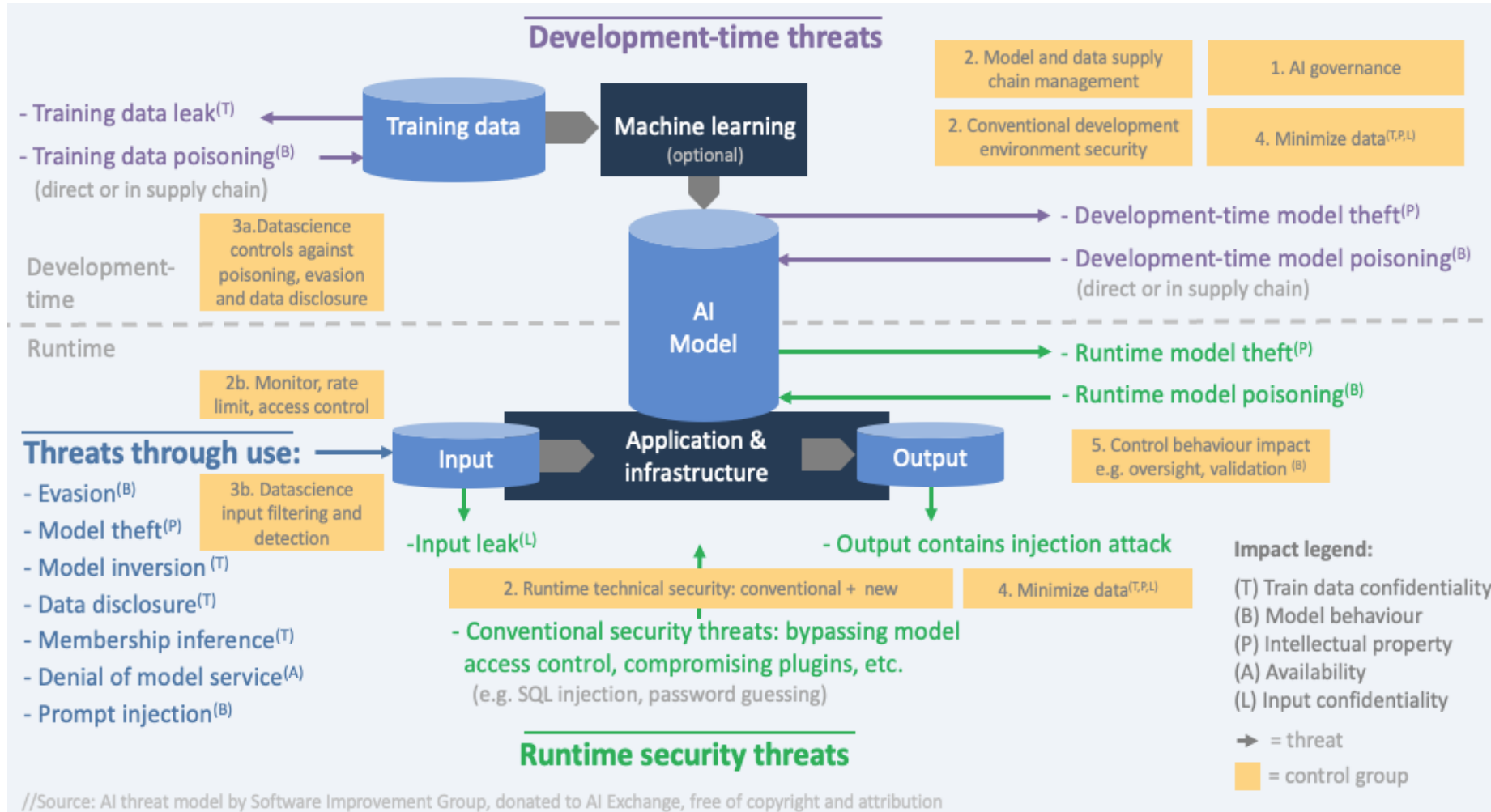
Impact on Confidentiality, Integrity or Availability



Threat and Impact



Threat and Impact with Controls



The periodic table of AI security

The periodic table of AI security



Found at: <https://owaspai.org/qoto/periodictable/>

The table below, created by the OWASP AI Exchange, shows the various threats to AI and the controls you can use against them – all organized by asset, impact and attack surface, with deeplinks to comprehensive coverage at the [AI Exchange website](#) with further references to related standards.

Note that [general governance controls](#) apply to all threats.

- The various threats to AI and the controls organized by asset, Impact and attack surface

Asset & Impact	Attack surface with lifecycle	Threat/Risk category	Controls	
Model behaviour Integrity	Runtime - Model use (provide input/ read output)	Direct prompt injection	Limit unwanted behavior , Input validation , further controls implemented in the model itself	
		Indirect prompt injection	Input validation , Input segregation	
		Evasion (e.g. adversarial examples)	Limit unwanted behavior , Monitor , rate limit , model access control plus: Detect odd input , detect adversarial input , evasion robust model , train adversarial , input distortion , adversarial robust distillation	
	Runtime - Break into deployed model	Model poisoning runtime (reprogramming)	Limit unwanted behavior , Runtime model integrity , runtime model input/output integrity	
	Development - Engineering		Model poisoning development time	Limit unwanted behavior , Development environment security , data segregation , federated learning , supply chain management plus: model ensemble
				Limit unwanted behavior , Development environment security , data



German
OWASP
Day 2024

The following threats and controls are highlights from the AI Exchange of which most are not in the LLM top 10.



German
OWASP
Day 2024

1

General Controls



Governance

Controls
AI Program
Security Program
Secure Development Program
Development Program
Check Compliance
Security Education



Limit the effects of unwanted behavior

Control
Oversight
Least Privilege
AI Transparency
Continuous Validation
Explainability
Unwanted Bias Testing



Sensitive Data Limitation

Control
Data Minimization
Allowed Data
Short Retain
Discrete



German
OWASP
Day 2024

2

Threats Through Use



Model Behavior Manipulation

Threat and Impact
Evasion
Closed-box evasion
Open-box evasion
Evasion After Data Poisoning



German
OWASP
Day 2024

3

Development-Time Threats



Control	Description
Development Security	Sufficient security of the AI development infrastructure, also taking into account the sensitive information that is typical to AI: training data, test data, model parameters and technical documentation
Segregate Data	Store sensitive development data (training or test data, model parameters, technical documentation) in a separated areas with restricted access.
Confidential Compute	If available and possible, use features of the data science execution environment to hide training data and model parameters from model engineers - even while it is in use.
Federated Learning	Federated learning can be applied when a training set is distributed over different organizations, preventing that the data needs to be collected in a central place - increasing the risk of leaking.
Supply Chain Management	Managing the supply chain to minimize the security risk from externally obtained elements. In regular software engineering these elements are source code or software components (e.g. open source).



Sensitive Data Leak During Development

Threat and Impact	Description
Development-time data leak	Unauthorized access to train or test data through a data leak of the development environment. This has an Impact on the confidentiality breach of sensitive train/test data.
Model theft through development-time model parameter leak	Unauthorized access to model parameters through a data leak of the development environment. This has an impact on the confidentiality breach of model intellectual property.
Source code/configuration leak	Unauthorized access to code or configuration that leads to the model, through a data leak of the development environment. Such code or configuration is used to preprocess the training/test data and train the model. This has a direct impact on confidentiality breach of model intellectual property.



German
OWASP
Day 2024

4

Runtime Application Security Threats



Leak Sensitive Input Data (*)

Control	Description
Leak Sensitive Input Data	<p>Input data can be sensitive (e.g. GenAI prompts) and can either leak through a failure or through an attack, such as a man-in-the-middle attack.</p> <p>Impact: Confidentiality breach of sensitive input data.</p>

Roadmap



Key Deliverables

- Prep 1.0: Review by community and by ourselves -> release 1.0
- Feed the Exchange 1.0 into at least the AI Act and ISO 27090
- Make it easier for readers to recognize their deployment model and select only what is relevant to them
- More illustration of threat models and attack vectors
- Further alignment with Mitre Atlas, NIST, the LLM Top 10, ENISA's work, and the AIAPP International Privacy Group



Get Involved and Contribute

Engage with the OWASP AI team through various platforms.

- **Connect** with us on the [OWASP Slack](#) workspace in the [#project-ai-community](#) channel. Authors are in the closed [#project-ai-authors](#) channel.
- Keep up with the latest **updates** by following us on [Twitter](#) and [LinkedIn](#).
- For technical inquiries and suggestions, **participate** in our [GitHub Discussions](#), or report and track issues on [GitHub Issues](#).

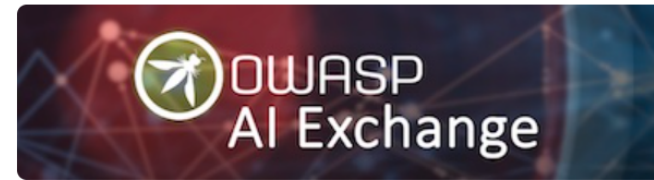
If contributing interests you, check out our [Contribution Guidelines](#) or get in touch with our project leaders.

The Exchange is built on expertise from contributors around the world and across all disciplines.



Where can I find more information?

OWASPAI.ORG



Comprehensive guidance and alignment on how to protect AI against security threats - by professionals, for professionals.

 Charter

 Connect with us!

 Contribute

 Register

 Media

 Navigator

Our Content

AI Security Overview

1. General controls

2. Threats through use

3. Development-time threats

4. Runtime application security threats



Participate in Content Development



- 📄 Send your suggestion to the [project leader](#).
- 🗋️ Join [#project-ai-community](#) in our [Slack](#) workspace.
- 🗣️ Discuss with the [project leader](#) how to become part of the writing group.
- 💡 Propose your [concepts](#), or submit an [issue](#).
- 📁 Fork our repo and submit a [Pull Request](#) for concrete fixes (e.g. grammar/typos) or content already approved by the core team.
- 🙌 Showcase your [contributions](#).
- 🐛 Identify an [issue](#) or fix it on a [Pull Request](#).
- 💬 Provide your insights in [GitHub Discussions](#).
- 🙋 Pose your [questions](#).



German
OWASP
Day 2024

OWASP AI Exchange

Follow us on LinkedIn





Stream the bi-weekly meetings





Get Involved as part of the Slack communication channel

(over 350 members)



project-ai-community



Get Notifications for @ Mentions



Huddle

About

Members 257

Integrations

Settings

Topic

Edit

Artificial Intelligence security and privacy

Description

Edit

Channel for the community interested in the content and the direction of the OWASP AI Exchange: <https://owaspai.org>

Managed by



Rob van der Veer(SIG)

Created by

Rob van der Veer(SIG) on December 20, 2022



FAQ

What is our the stance on privacy?

https://owaspai.org/docs/ai_security_overview/#how-about-privacy

What is our the stance on copyright?

https://owaspai.org/docs/ai_security_overview/#how-about-copyright

How can I get associated?

<https://owaspai.org/contribute/>

What are all the frameworks that exist around GenAI?

<https://owaspai.org/goto/references/>

How does it complement OWASP LLM Top Ten?

The LLM is about the top 10 issues. **The Exchange is about all issues in all of AI**





Key Reference Links

- [Bi-Weekly Meeting](#)
- [Contribute](#)
- [OWASP Slack Invite](#)
- [OWASP LLM top 10](#)
- [ENISA ML threats and countermeasures 2021](#)
- [MITRE ATLAS framework for AI threats](#)
- [NIST threat taxonomy](#)
- [ETSI SAI Problem statement Section 6](#)
- [Microsoft AI failure modes](#)
- [NIST](#)
- [NISTIR 8269 - A Taxonomy and Terminology of Adversarial Machine Learning](#)
- [OWASP ML top 10](#)
- [PLOT4ai threat library](#)
- [AVID AI Vulnerability database](#)
- [OECD AI Incidents Monitor \(AIM\)](#)
- [ENISA AI security standard discussion](#)
- [ENISA's multilayer AI security framework](#)
- [Alan Turing institute's AI standards hub](#)
- [Microsoft/MITRE tooling for ML teams](#)
- [Google's Secure AI Framework](#)
- [NIST AI Risk Management Framework 1.0](#)
- [ETSI GR SAI 002 V 1.1.1 Securing Artificial Intelligence \(SAI\) – Data Supply Chain Security](#)
- [ISO/IEC 20547-4 Big data security](#)
- [IEEE 2813 Big Data Business Security Risk Assessment](#)
- [BIML](#)
- [Media](#)
- [OWASPAI.ORG](#)

A futuristic, metallic bee with a black and gold striped abdomen and a highly detailed, mechanical head and legs. The bee is standing on a glowing, reflective digital platform. The background features a cityscape at sunset with a warm orange glow, and various glowing icons and symbols, including a spider, a heart, and a Wi-Fi symbol, floating in the air. The overall scene is a blend of nature and technology.

OWASP AI Exchange
(owaspai.org)

Thank You