Feedback from a year of work with our clients

# How can cyber help make AI a success?

Thomas ARGHERIA
AI Security Manager
thomas.argheria@wavestone.com
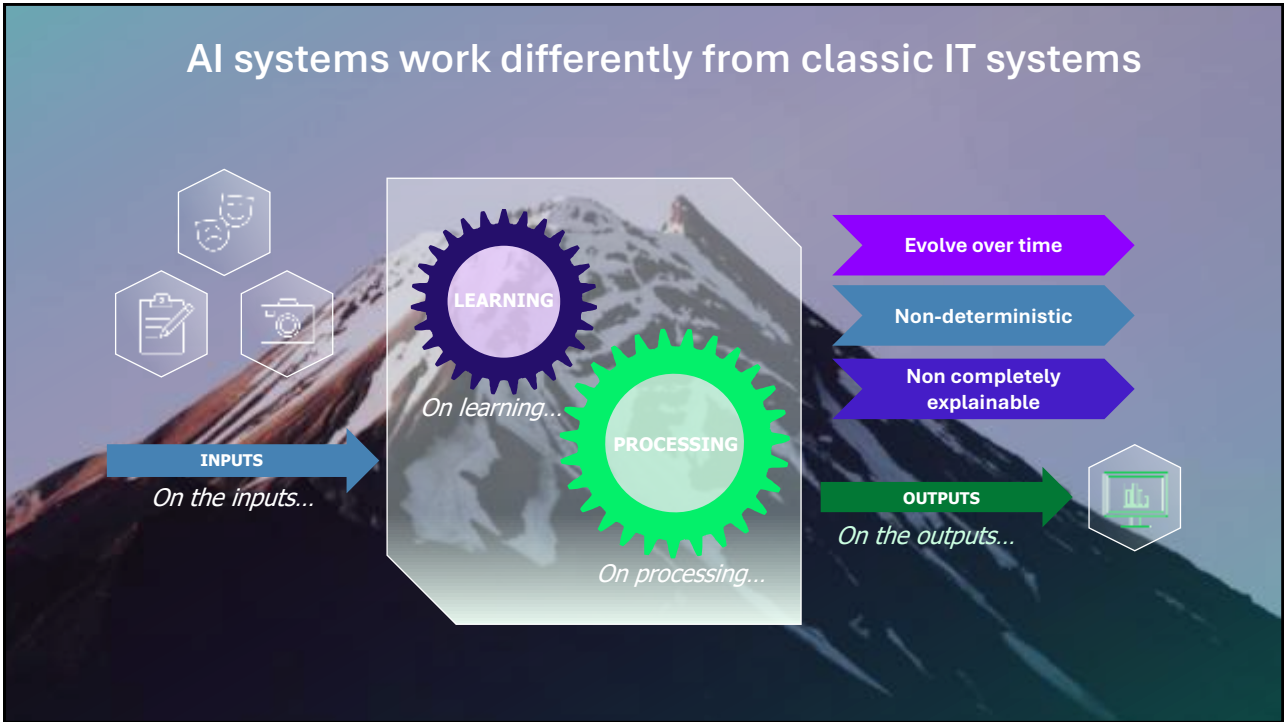
1



# No doubt: AI is a unique opportunity...

## ...that must be secured!

2

# AI systems work differently from classic IT systems



# AI can be attacked in very specific and new ways...

# New regulations arise and include cybersecurity expectations

| UNITED-STATES | EUROPE | CHINA |
|---|---|---|
| **Executive Order** | **AI ACT** | **Cybersecurity requirements for GenAI services** |
| *In place since October 2023* | *In place since March 2024* | *In place since May 2024* |
| An approach focused **on market self-regulation** | The EU positioned itself as the world's police officer and **push for citizen protection** | China **focus on pushing for best practices** in AI management and data management |
| • **Light** AI security regulation<br>• Focus on **guidelines** for administrations<br>• Let the states define their own approach | • **Risk-based** approach<br>• Every organization must comply by **May 2027**.<br>• **Already some consequences**: new iPhone with GenAI & ChatGPT voice chat functionality postponed ... | • China is focusing on the **cybersecurity of its system with a risk-based approach** and on **regulating the processing of data, especially labeling** |

5

---

# Today, though the hype effect, AI is a reality!

Some clients are adopting AI on a large scale:

• **Between 50 and 400** uses cases identified

• A strong **mobilization at Excom level**

Leading to a **lot of activities** but a **lot of blurriness.**

## Our goal: help clarify how to tackle the AI security topic

Worked with **+20 clients already working on the topic.**

We **benched our clients on their AI maturity,** based on the 5 NIST's pillars

• Govern
• Identify
• Protect
• Detect
• Respond

© WAVESTONE | 6

6

## First lesson: state your stance on AI!

**AI Advanced Creators** — **35%** of our clients

- **Build and sometimes sell AI models**
- Both **third party and in-house solutions**
- **Structured teams of data scientists** and proven data science processes.

**AI Orchestrators** — **35%** of our clients

- Embeds **AI functionalities** in their products/services, internally or externally.
- **Make available a GenAI Platform** for app builders
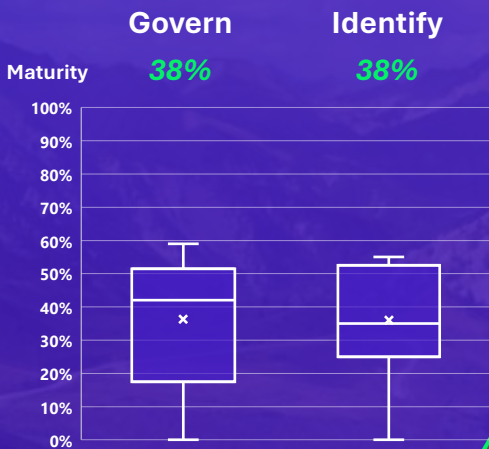- Mostly use **third party solutions**, that they integrate

**AI Users** — **30%** of our clients

- Uses AI **punctually to boost productivity**
- **Uses third-party solutions**
- No structured teams of data scientist or AI Hub.

7

## *Market quickly embraced the need to adapt for AI's arrival*

| | Govern | Identify |
|---|---|---|
| Maturity | *38%* | *38%* |

*Source: Wavestone AI CyberBenchmark 2024*

WAVESTONE

8

# A new governance to define at group level, with few resources
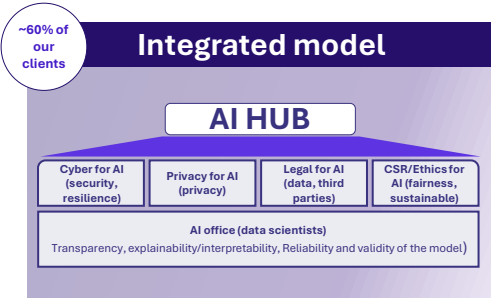
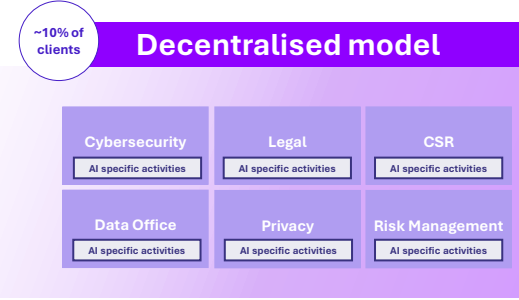**86%** Have a defined **trustworthy governance at group level**

**14%** Of companies have **sufficient AI expertise in regards with the stakes**

Our recommendation: **compensate with an integrated governance** that will help people augment their skills

**~60% of our clients**

## Integrated model

**~10% of clients**

## Decentralised model

**AI HUB**

| Cyber for AI (security, resilience) | Privacy for AI (privacy) | Legal for AI (data, third parties) | CSR/Ethics for AI (fairness, sustainable) |
|---|---|---|---|

**AI office (data scientists)**
Transparency, explainability/interpretability, Reliability and validity of the model)

**~30% of clients in hybrid mode**

| Cybersecurity | Legal | CSR |
|---|---|---|
| AI specific activities | AI specific activities | AI specific activities |
| Data Office | Privacy | Risk Management |
| AI specific activities | AI specific activities | AI specific activities |

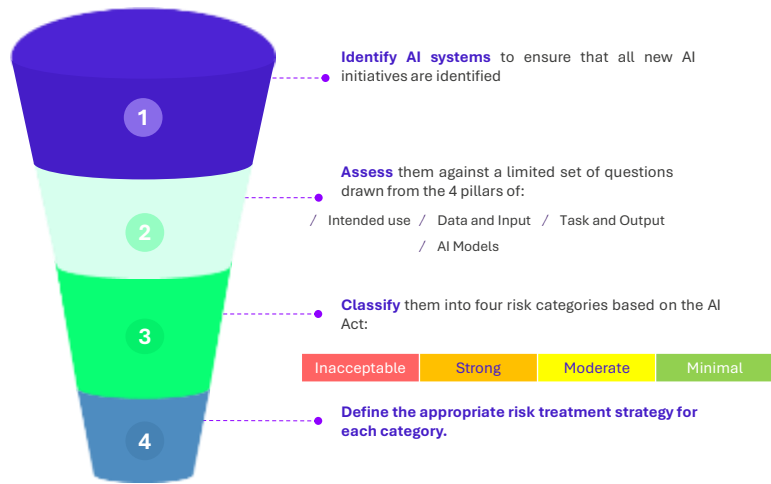© WAVESTONE | 9

9

# Frame your cyber approach

**64%** Have an **AI security policy**
- *Frame use of AI large public application*
- *Indicates the process to secure AI project*
- *Integrate Thrid Party stance against AI*

**71%** Have adapted their project **processes** for AI
- *Define role and responsibilities*
- *Define validation process*

**1** **Identify AI systems** to ensure that all new AI initiatives are identified

**2** **Assess** them against a limited set of questions drawn from the 4 pillars of:
/ Intended use / Data and Input / Task and Output
/ AI Models

**3** **Classify** them into four risk categories based on the AI Act:

| Inacceptable | Strong | Moderate | Minimal |
|---|---|---|---|

**4** **Define the appropriate risk treatment strategy for each category.**

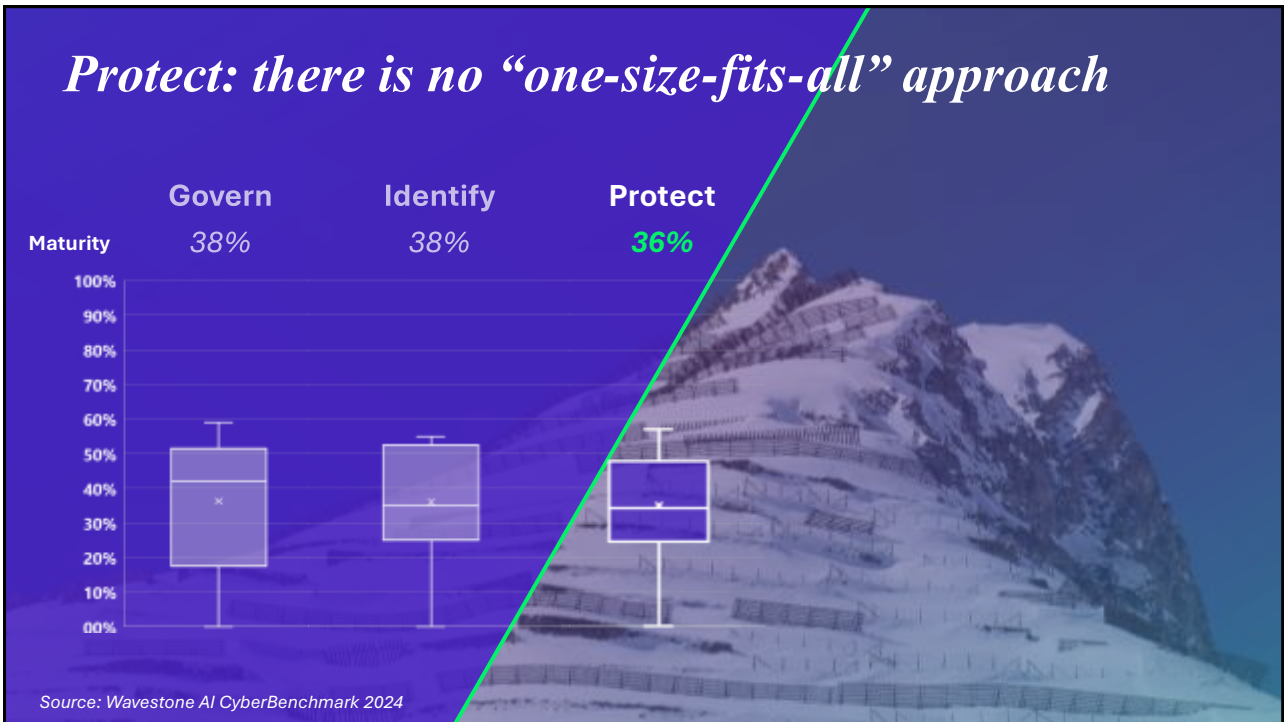**Wavestone Accelerators** — Assessment questionnaire — Risk level analysis

10

## We identified the key recurring factors responsible for the greatest risks

**Six common red flags**

- External facing systems, especially GenAI chatbot
- Dataset for training unknown or containing personal data
- Retrieval Augmented Generation (RAG) on critical / confidential data
- Model modifications, sources or toolset from non-authoritative sources
- GenAI capability to take actions
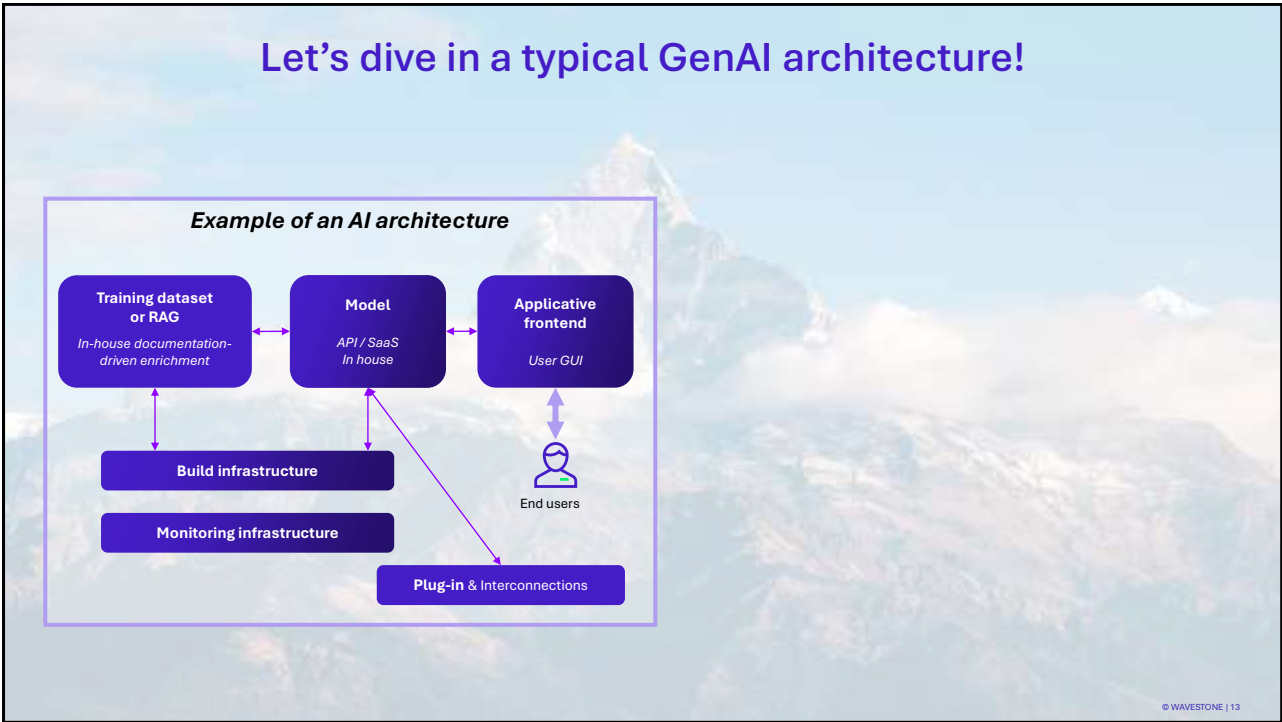- AI model with mission critical output (safety detection for instance)

**But most of AI** use case we assessed are typically used for **non-critical processes** that don't demand high availability or strict integrity, often relying on human oversight

11

## *Protect: there is no "one-size-fits-all" approach*

| Govern | Identify | Protect |
|--------|----------|---------|
| *38%* | *38%* | ***36%*** |

Maturity

*Source: Wavestone AI CyberBenchmark 2024*

12

## Let's dive in a typical GenAI architecture!



*Example of an AI architecture*

13

## AI users: secure your data and check your suppliers



*Example of an AI architecture*

- **Protect the data** being accessed or generated (access rights, policies, etc.)
- **Configure the parameters** and ensure the ability to monitor the ecosystem
- **Select your providers**: verify compliance with your security requirements (learning phase, data usage, etc.) **including contractual** requirements and measures regarding shared data

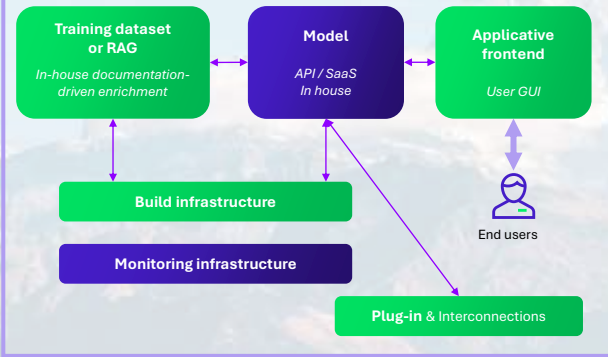**43%** 43% of clients **adapted their Third Party assessment** methodology for AI vendors

Component to protect

14

## AI Orchestrator: choose your models and platforms and implement MLSecOps

### Example of an AI architecture

Training dataset or RAG
*In-house documentation-driven enrichment*

Model
*API / SaaS*
*In house*

Applicative frontend
*User GUI*

Build infrastructure

Monitoring infrastructure

End users

Plug-in & Interconnections

Component to protect

- **Set up criteria** to choose the right model: whitelist suppliers, code review, operational testing...
- Build **inputs and output controls**
- Ensure proper **security of the front end**
- Make AI project "**secure by design**" with MLSecOps

**42%** Have a **model selection process to identify trusted sources**

© WAVESTONE | 15

15

## Advanced Creators: full responsibility of the whole stack

### Example of an AI architecture

Training dataset or RAG
*In-house documentation-driven enrichment*

Model
*API / SaaS*
*In house*

Applicative frontend
*User GUI*

Build infrastructure

Monitoring infrastructure

End users
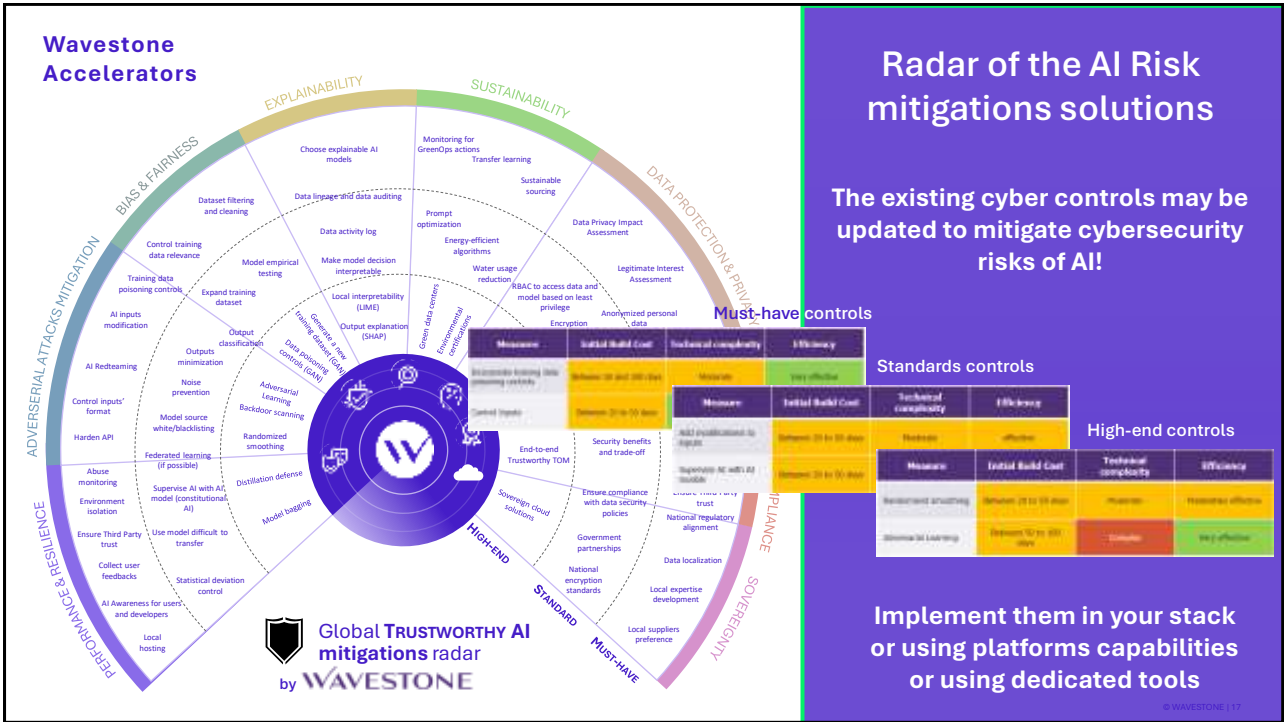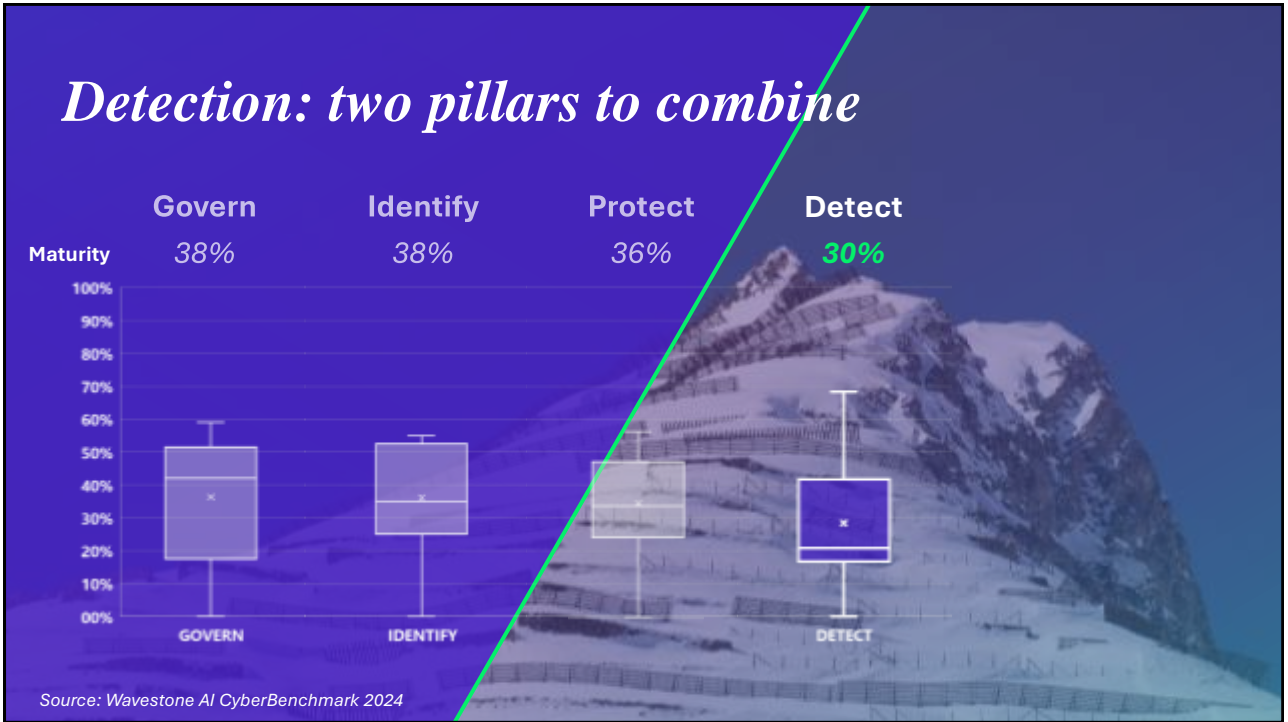
Plug-in & Interconnections

Component to protect

- Implement **in-depth security measures**, alongside with data scientist:
  - **Model architecture security**, as randomized smoothing, adversarial learning, bagging
  - **Training data security**: with synthetic data, differential privacy
  - **Model protection**, as homomorphic encryption and differential privacy...
- Think about the security measures as a **differentiator** to resell your apps and model

**7%** Have established measures and adapted **tooling to detect and defend** against **malicious prompts** and other **identified threats**

© WAVESTONE | 16

16

17



18

## *Detection: two pillars to combine*

| Maturity | Govern | Identify | Protect | Detect |
|---|---|---|---|---|
| | *38%* | *38%* | *36%* | *30%* |

*Source: Wavestone AI CyberBenchmark 2024*

19

## First, Pentesting! But with a twist: threats are present along the entire AI lifecycle

Collection → Processing → Model → Tests → Deployment → Monitoring

**Poisoning attacks**
/ Dataset poisoning
/ Retraining poisoning

**Oracle attacks**
/ Membership inference
/ Model extraction
/ Model inversion

**Manipulation attacks**
/ Evasion
/ Model reprogramming
/ Denial of service

**Prompt injection**

... that we tested and adapted to land **our AI redteam framework** on the market

| Assessing AI capabilities and biases | Assessing AI systems flaws |
|---|---|
| *Hallucination, Misinformation, Robustness, Harmfulness Prompt Injection...* | *Pre-prompt access, Input/Output filtering, Illegitimate internal data retrieval , API limitations, Detection & monitoring* |

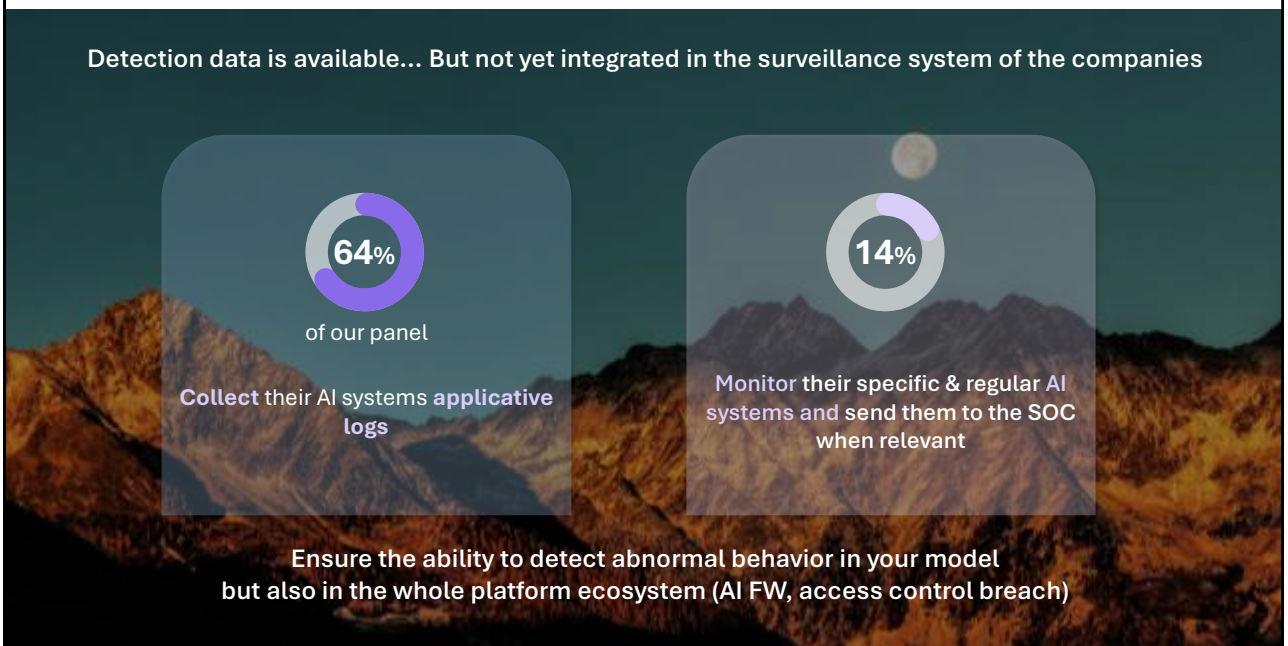**New approach and tooling required, often using LLM to attack LLM!**

**64%** have a pentest process in place to test the use case

**7%** Use advanced model robustness evaluation

© WAVESTONE | 20

20

## Feedback from our GenAI Red Teaming team

### TOP FLAWS IDENTIFIED

**+10 PROJECTS**
Chatbots, GenAI, LLM, etc.

**7 SECTORS**
Energy, retail, luxury, transportation, chemicals, cosmetics, distribution.

**100% JAILBROKEN**
Illegitimate content, hallucination, bias, etc.

Web Integration flaws & Injection attacks

Weak privileges management

Lack of monitoring

Faulty DevSecOps processes

**Standard**

Data leakage through prompt injection / trapped documents

ML/AI platform missing security configuration

API/Plugin security gaps

Overreliance on platform moderation

**ML Specific**

21

## Then, integrate AI systems in the global detection strategy

Detection data is available... But not yet integrated in the surveillance system of the companies

**64%**
of our panel

**Collect** their AI systems **applicative logs**

**14%**

**Monitor** their specific & regular AI systems and send them to the SOC when relevant

Ensure the ability to detect abnormal behavior in your model
but also in the whole platform ecosystem (AI FW, access control breach)

22

11

23



24

Source: Wavestone AI CyberBenchmark 2024

25



26

# One More Thing…

# AI can also enhance cybersecurity capabilities

**WAVESTONE**

27

## In short, there are 3 categories of use-cases to remember

**Ease communication activities**
1. Multi-language awareness
2. CISO GPT Chatbot to ease documentation access

**Accelerate processes**
3. Third Party Security questionnaires analysis
4. Data classification recommendation

**Reinvent detection and reaction**
5. GenAI SoC Copilots
6. SOC playbook update via ML

### Use Case Analysis Matrix

High · ADDED VALUE · Low

Difficult · *FEASABILITY* · Easy

*Use case highlighted are offered by at least one software vendor*

28

14

## A strong experience feedback on AI Security

## ... leveraging on thoughtful accelerators

AI STRATEGY & ROADMAP DEFINITION

AI RED-TEAMING

AI SECURITY POLICIES

AI PLATFORM SECURE IMPLEMENTATION

DATA GOVERNANCE DESIGN & ROLL-OUT

INTEGRATION OF SECURITY IN AI PROJECT

DESIGN & DEV OF CYBER AI USE CASE

Scan me to get the slides or join the benchmark

Meet us at booth **L31**!

**Deepfake me!**
*For awareness campaign on new deceptive social engineering campaign*

**HackMyAI**
*Facial Recognition hacking and RAG poisoning demonstrator*

**CISO GPT**
*To facilitate access to cybersecurity knowledge*

**Crisis Maker (alpha)**
*AI-enhanced crisis exercise*

29



WAVESTONE

30