# What is your AI risk Score?

# Legal notes and disclaimer

# The Foundation for a Successful AI Business Strategy - Trusthworthy AI



Mike Kehoe,
*EMEA Data Security Leader, Data & AI Platform*



Aleksandra Kaszuba
*NCEE Wx.Governance Virtual SWT leader*



Nicola Piazza
*EMEA Principal Security Technical Leader*



Albert Puah
*WW Sales Leader, Data Security*

# AI Governance and AI Security

The Sandwich Dilemma



Undesired Results

Desired Results

AI Governance

AI Security

Human  (Task Performing )

Ai (Task Performing )

https://food.ec.europa.eu/horizontal-topics/farm-fork-strategy_en

# The Data & AI Platform architecture

**Data Sources**

**Data Pipelines**

**Data Integration**

**Data Pipelines**

**Data Warehouses**

**Data Intelligence**

**Insights Apps**

IBM Cognos Analytics

## Data Sources

- Source System
- Source System
- IBM DB2™
- salesforce
- ORACLE

## Data Integration

- E**T**L Datastage
- ELT Datastage
  - DBT
- IBM DataStage®
- Data Virtualization
- Real time replication
  - Kafka
  - DB
- Real Time Streaming

IBM Streamsets

## Data Warehouses

- EDW BS
- Datalake OS
- EDW

**Lakehouse**

**watsonx.data**

IBM watsonx.data™

- Poli-engine
  Presto, Spark, DB2, NZ, DataStax
- Metadata Exchange
  WKC connector & other
- Data Object Store Iceberg
  Avro, Parquet, Orc

## Data Intelligence

### Knowledge Catalog

- Data Discovery
- Data Privacy rules
- Tags, Business Glossary
- Data Quality rules
- Lineage
- Master Data (Agile)

Knowledge Catalogs

## Insights Apps

- Business Intelligence Cognos Analytics
- Planning Analytics
- Power BI

watson**x** Code Assistant
watson**x** Assistant
watson**x** Orchestrate
watson**x** Orders

IBM watsonx Code Assistant™

IBM watsonx Orchestrate™

IBM Business Automation

IBM watsonx.ai™

## GIGO

## Great In Garbage Out

**AI Platform**

**watsonx.ai**

### Traditional AI / ML
Auto AI
Model Serving
Data Science
Model Ops
Optimization

### Generative AI – Foundational Models
Zero | Few | Many Shot
Prompt Lab
Prompt Optimization
Model tuning

**AI Governance**

**watsonx.governance**

IBM watsonx.governance™

- Model Inventory & Performance
- Model Risk Management
- Model evaluation & monitoring

Azure Machine Learning

Amazon SageMaker

IBM Guardium® Data Security Center

**Data Observability**

IBM Databand™

**Data Security**

Regulatory Fines

| Noncompliance case | Proposed fine |
|---|---|
| Breach of AI Act prohibitions | Fines up to €35 million or 7% of total worldwide annual turnover (revenue), whichever is higher |
| Noncompliance with the obligations set out for providers of high-risk AI systems or GPAI models, authorized representatives, importers, distributors, users or notified bodies | Fines up to €15 million or 3% of total worldwide annual turnover (revenue), whichever is higher |
| Supply of incorrect or misleading information to the notified bodies or national competent authorities in reply to a request | Fines up to €7.5 million or 1.5% of total worldwide annual turnover (revenue), whichever is higher |

https://artificialintelligenceact.eu/

# The reasons we MUST have Trustworthy AI

Regulatory Fines

Reputational Damage



DPD error caused chatbot to swear at customer

19 January 2024

Tom Gerken Technology reporter

Share    Save

The customer then posted the chat, which had gone viral with 1.3 million views and over 20 thousand likes.

# The reasons we MUST have Trustworthy AI

Regulatory Fines
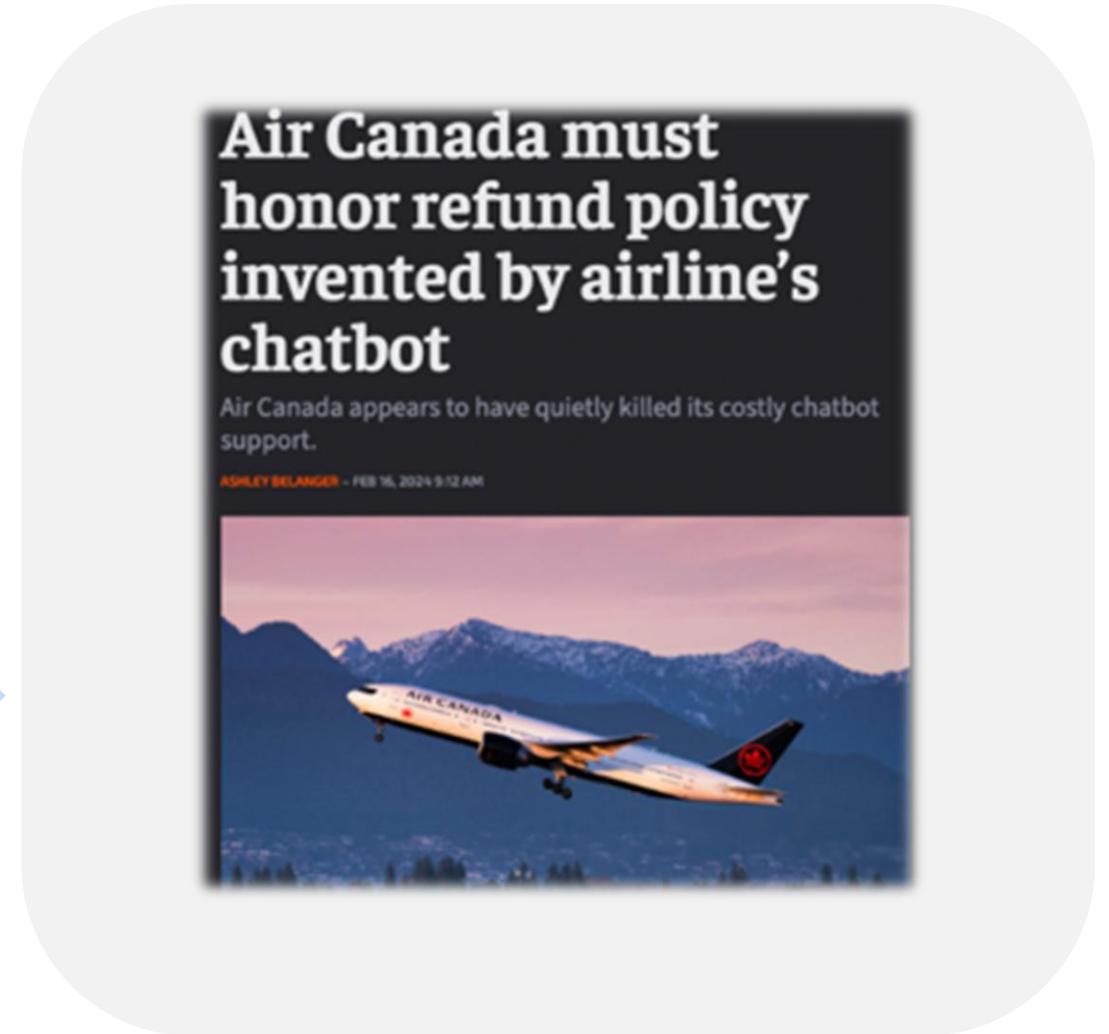
Reputational Damage

Revenue Loss
( primary / secondary )



Air Canada must honor refund policy invented by airline's chatbot

Air Canada appears to have quietly killed its costly chatbot support.

ASHLEY BELANGER – FEB 16, 2024 9:12 AM

# The reasons we MUST have Trustworthy AI                    The 4 Rs

Regulatory Fines

Reputational Damage

Revenue Loss
*( primary / secondary )*

Running Operations Impacts
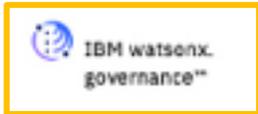
**Gartner**  Who We Serve  Our Solutions  Latest Insight  Conferences  Webinars  AI

Newsroom  Topics  Media Contacts  Media Resources  Insights  Archive

Newsroom / Information Technology / Press Release

## Gartner Predicts Over 40% of Agentic AI Projects Will Be Canceled by End of 2027

SYDNEY, Australia, June 25, 2025

**Analysts to Explore Agentic AI Trends During Gartner IT Symposium/Xpo, September 8-10 on the Gold Coast**

due to:-

- escalating costs
- unclear business value
- inadequate risk controls

https://www.gartner.com/en/newsroom/press-releases/2025-06-25-gartner-predicts-over-40-percent-of-agentic-ai-projects-will-be-canceled-by-end-of-2027

# Evolution of AI and its vulnerabilities



Here's a drawing of an IBM pickup truck in Dublin:
Do you need any other images?

# Evolution of AI and its vulnerabilities



Accountability

Accuracy

Fairness

Veracity

Transparency

Drift

Trusted data

Energy consumption

Explainability

Adversarial Robustness

IP/PII leakage

...

IBM watsonx. governance™

**Exploit difficulty**

**Model extraction**
Steal a model's behavior by observing the relationships between inputs and outputs

**Inversion exploits**
Reveal information on the data used to train a model, despite only having access to the model itself

**Backdoor exploits**
Alter a model subtly during training to cause unintended behaviors under certain triggers

**Data poisoning**
Change the behavior of AI models by altering the data used to train them

**Model evasion**
Circumvent the intended behavior of an AI model by crafting inputs that trick it

**Supply chain exploits**
Generate harmful models that hide malicious behavior, or target vulnerabilities in systems connected to the AI models

**Prompt injection**
Manipulate AI models into performing unintended actions by dropping guardrails and limitations put in place by the developers

**Data exfiltration**
Access and steal sensitive data used in training and tuning models through vulnerabilities, phishing, or misused privilege credentials

**Potential impact**

IBM Guardium® Data Security Center

# How AI can misbehave

*(these are most common attributes when AI misbehaves )*

## Passive Anomaly
( accidental )

- ## Quality Assurance Miss

  - Clumsy & Lossy

  - No Malicious Intent

  - Remediation needs no escalation

  - Root cause Analysis

  - Maintainable Impacts

## Active Anomaly
(coerced)

- ## Targeted system abuse

  - Focused & Covert

  - Malicious Intent

  - Remediation needs escalation due to deeper security implications

  - Extensive Root cause Analysis needed

  - Significant Impact

**AI Governance (in->out)**

**AI Security (in<-out)**

**Trustworthy AI**

# Better Together

## Guardium AI Security + watsonx.governance – Integrated value

| Capability | Guardium AI Security | Watsonx.gov | Integrated Value |
|---|---|---|---|
| AI Discovery | Automatically detects Shadow AI and inventories ML/LLM/foundation models/Ai Software/AI Services/AI Agents and more | Tracks and manages known and registered AI assets in the model catalog. | **Unified, automated and continuously up-to-date AI catalog.** Complete visibility into all AI systems and assured governance without relying on human discretion. |
| Risk & Compliance | Scans for AI policy & compliance violations and unauthorized behaviors and automates evidence collection. | Maps policies to regulatory frameworks (e.g., NIST, EU AI Act) | **Proactive risk and compliance management** – Continuous compliance validation and streamlined audit readiness. |
| AI Posture Management | Security Posture Management (SPM) - Scans for data, model, and usage vulnerabilities, misconfigurations, threat intel and provides remediation. | Responsible AI Posture Management - Monitors for bias, drift, fairness, and performance. | **Secure & Responsible AI -** Address ALL technical and ethical risks for AI system for all stakeholders involves with aspects of risk management. |
| Lifecycle Management | Pen-tests models, LLMs and AI systems and analyzes risk for safe third-party embedded AI onboarding. | Governs model lifecycle through approvals and version control. | **Secure and Governed AI Lifecycle Management -** Actively reduce risk and elevate safety in all AI deployments. |
| Operational Integration | Runtime Controls - runtime protection and AI access control (e.g., prompt injection attacks, PII protection and access to LLMs) | Facilitates collaboration across product, risk, and compliance teams. | **Integrated Operations -** Operationalize safety and security across product, risk, compliance, and security stakeholders. |
| AI Policy Management & Enforcement | End-to-end security policy management and enforcement through all stages of the lifecycle (development through production). | Centralized organizational policy management and controls for governing the AI use case and AI lifecycle. | **Policy definition, enforcement, and monitoring** for all stakeholders involved with AI systems in their respective "language" (Security / compliance), unified to complete a complete AI lifecycle policy. |

# Let's look at daily example – How it works in pratice



**2** Question: How can I join AMLA?

**watsonx.data**

**3** Chunks: best fitting fragments of documents that can answer this question

## Recruitment Assistant - AMLA

Application answers questions about AMLA's recruitment process using Azure OpenAI GPT :

Data is logged in watsonx.governance.

Enter your question:

How can I join AMLA?

Get Answer

### Response

Read the vacancy notice carefully, apply for jobs that match your experience and aspirations, ensure you meet all minimum requirements, and submit your application electronically before the deadline indicated in the advertised vacancy notice.

**1** How can I join AMLA?

**John Smith**
Candidate

**7** Answer: Read the vacancy notice carefully and …

**Guardium AI Security**

**4** Question + Chunks

**GPT-5**

**5** Answer: Read the vacancy notice carefully and apply for jobs that match your experience …

**6** Question + Chunks + Answer

**watsonx.governance**

For onboarding an external AI model, such as GPT-5, it is enough to depend only on the vendor's documentation.

IBM.

# False

## vendor docs matter, but the organization must adapt and govern use

Risk mapping in AI governance helps organizations connect compliance requirements with potential technical and ethical risks.

IBM®

**True**

# An AI risk score helps organizations evaluate both ethical and technical risks in AI systems.

Technical documentation is more than an AI Act requirement – it is essential for system transparency, risk assessment, and safe operations.

IBM.

# True it supports explainability, accountability, and safe AI governance.

# The General-Purpose AI (GPAI)
## Code of Practice

- The 3 chapters of the code - **Transparency, Copyright,** and **Safety** and **Security**.

- The Chapters on Transparency and Copyright offer all providers of general-purpose AI models a way to demonstrate compliance with their obligations under Article 53 AI Act.

- The Chapters on Safety and Security is only relevant to the small number of providers of the most advanced models, those that are subject to the AI Act's obligations for providers of general-purpose AI models with systemic risk under Article 55 AI Act.

**Signatories of Code of Practice (*):**

Accexible, AI Alignment Solutions, Aleph Alpha, Almawave, Amazon, Anthropic, Bria AI, Cohere, Cyber Institute, Domyn, Dweve, Euc Inovação Portugal, Fastweb, Google, Humane Technology, IBM, Lawise, Microsoft, Mistral AI, Open Hippo, OpenAI, Pleias, re-inventa, ServiceNow, Virtuo Turing, WRITER, In addition, xAI signed up to the Safety and Security Chapter; this means that it will have to demonstrate compliance with the AI Act's obligations concerning transparency and copyright via alternative adequate means.

European Commission

*\* Some signatories may not appear immediately, but we are making sure to continuously update the list as signatures are confirmed.*

# IBM AI Security full end to end Solution's features

- Automated responses
- Continuous monitoring of all AI assets
- Risk posture status
- Risk Management Frameworks adherence

Configuration

Inventory

Discovering

Scoring

Remediation

Monitoring & Respond

Posture Mgt

Pen testing

Compliance Reporting

# Take Away Points

New risk but the same flaws & threat actors

The need is the same , >>>> mitigate against the 4 Rs

AI for business strategy can't succeed without AI governance + AI security ( Trustworthy AI )

AI will generate many new opportunities for commercial success but  also new Malicious opportunities

Start with the end in mind …. design Ai + Governance + Ai Security as a single project

Seems overwhelming ?  No problem IBM Can help