

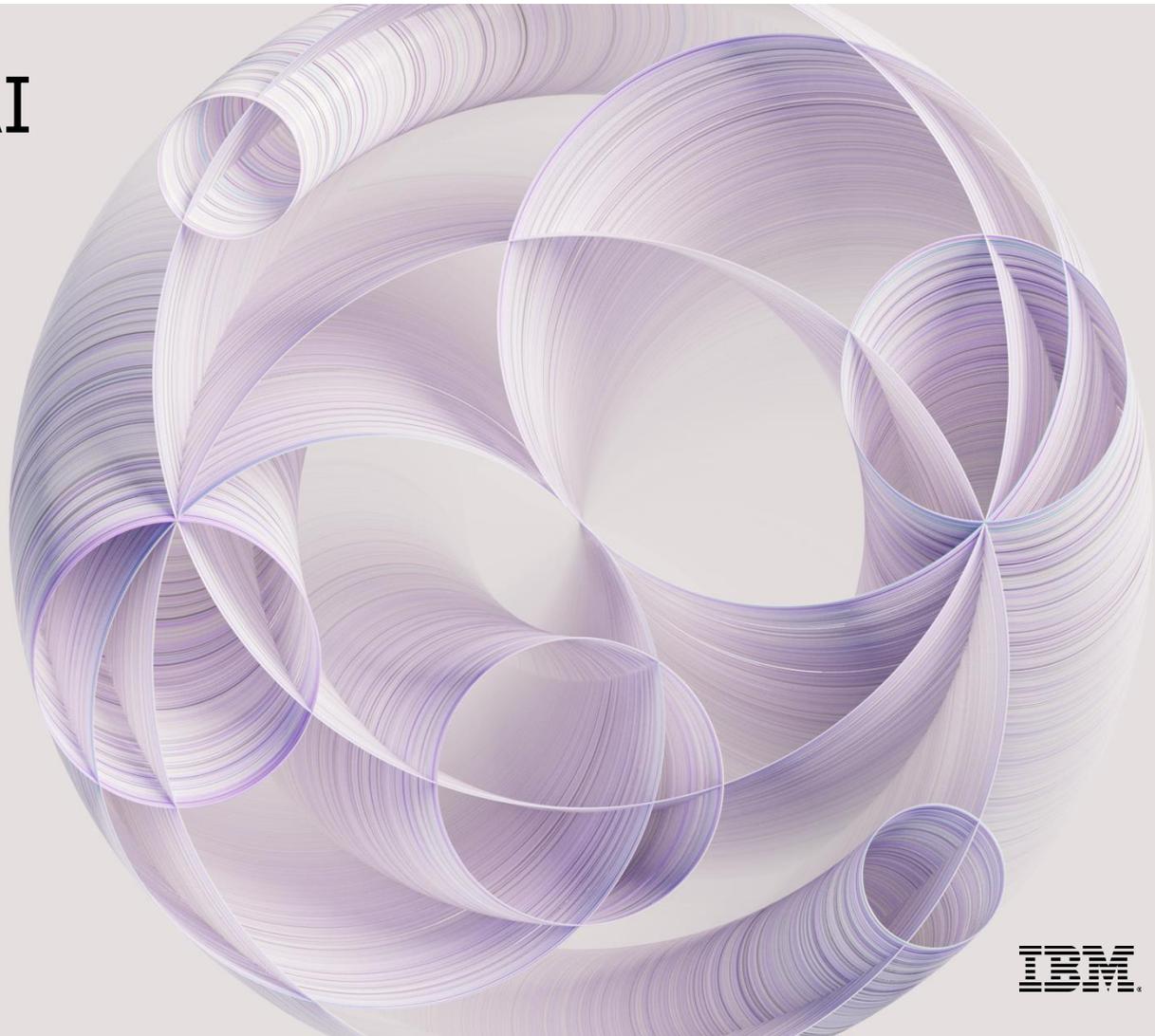
# Operationalizing AI

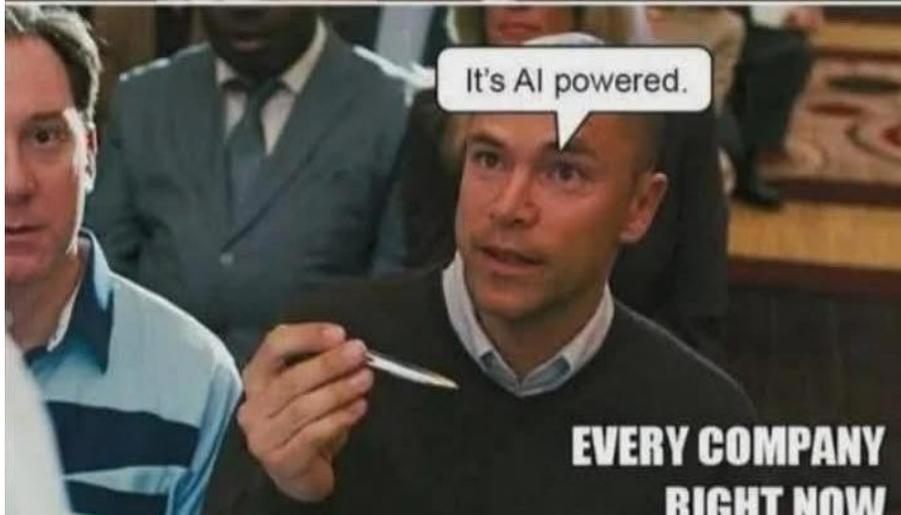


Hans-Petter Dalen  
Business Executive for EMEA

AI for Business

watsonx.

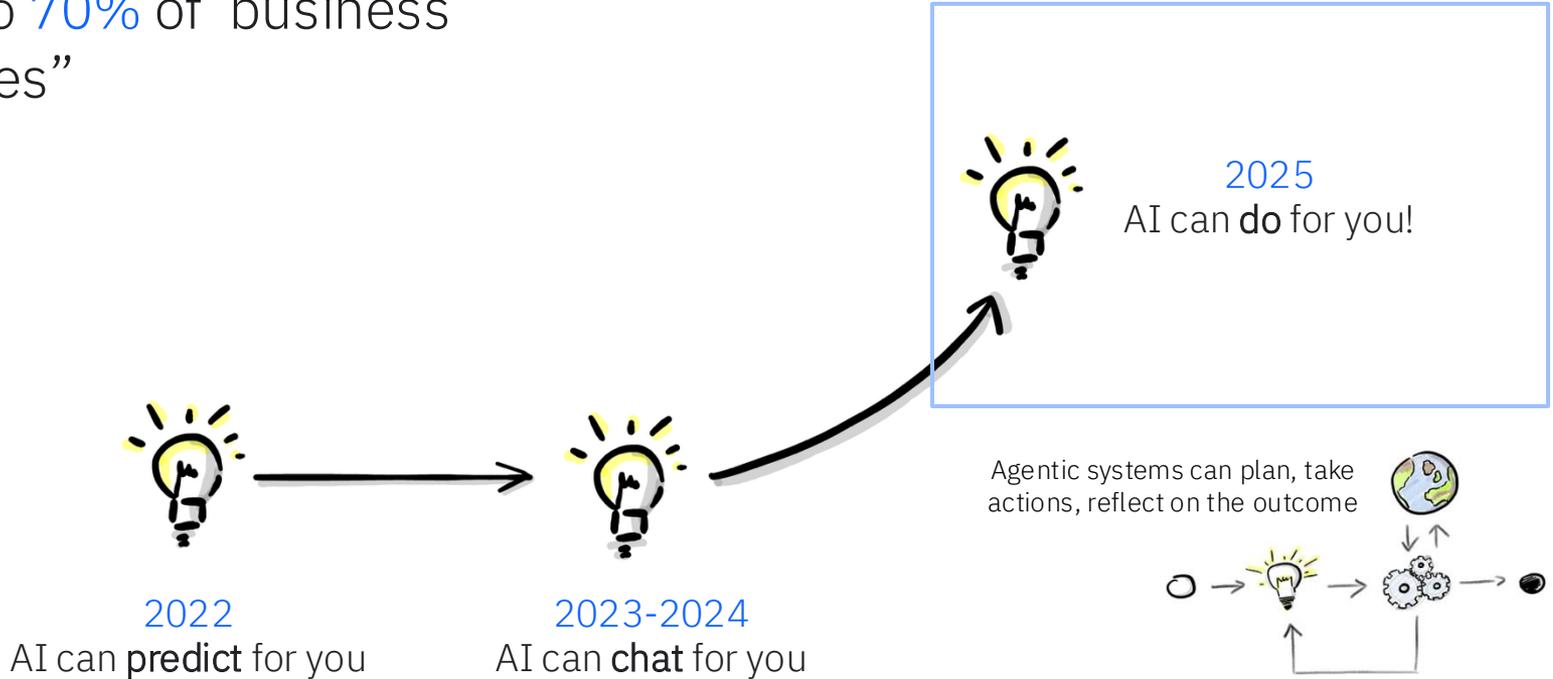




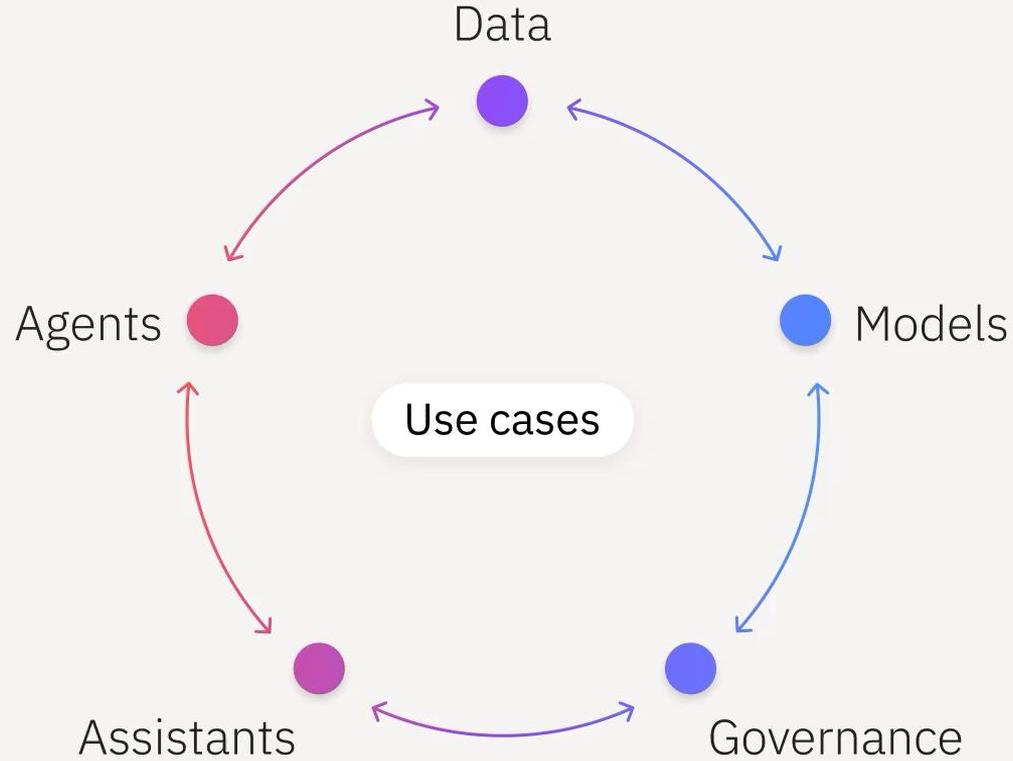
# Business impact of AI

“AI could enable automation of up to 70% of business activities”

(McKinsey)



# AI building blocks to the future



# IBM is “client zero” for AI governance at scale

## Results

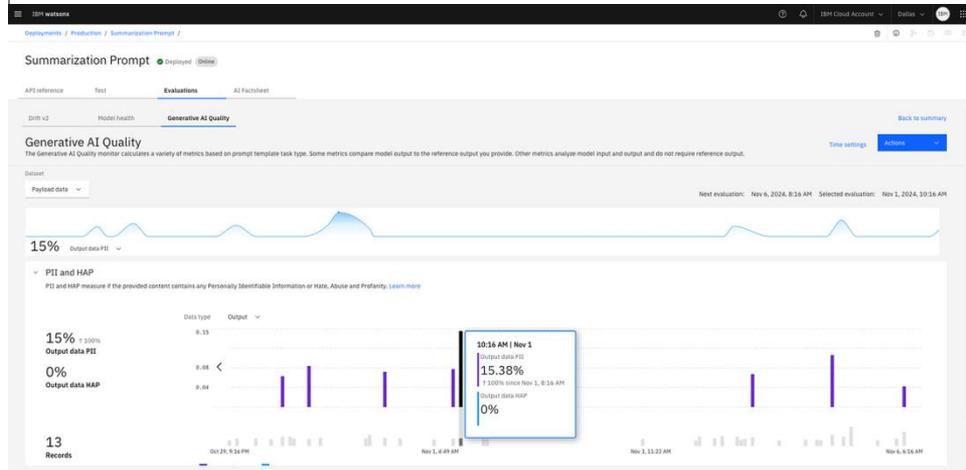
- >6,500 applications and processes managed
- **6 weeks** time for a recent new enterprise-wide compliance campaign
- **Days to minutes** projected reduction in time to evaluate AI solutions during build stage
- **2x** increase in Gen AI metrics evaluated during solution validation
- **14** health, drift, safety metrics continuously monitored in production
- Foundation for quick adoption of new technical advances

## People & Process

- Accountability model
- Centralized governance and decision making (AI Ethics Board)
- Focals in the business units
- Policies and Principles
- Ethics by Design
- Awareness and education
- Integrated Governance Program

## Technology

- watsonx.governance
- Integration with systems of record
- Self-service data gathering
- Standardized assessments
- Use case review process
- Extensive dashboards
- AI monitoring



Your principles of trust are the foundation of all your ethical decisions, including use case reviews

## IBM Principles for Trust and Transparency

- 1 The purpose of AI is to augment — not replace — human intelligence
- 2 Data and insights belong to their creator
- 3 New technology, including AI systems, must be transparent and explainable

## IBM Pillars of Trust



Explainability



Transparency



Fairness

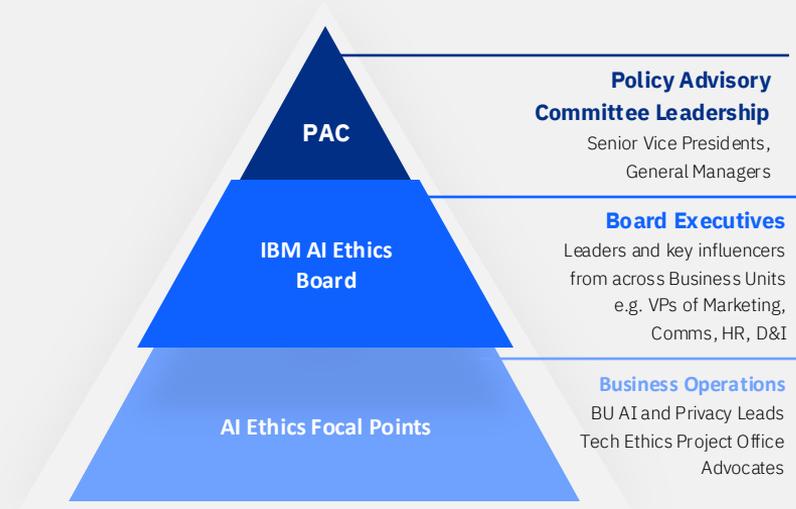


Privacy



Robustness

Centralized organizational AI governance enables top-down and bottom-up engagement



### Benefits of a centralized AI governance structure:

- Defined accountability and ownership starting at the top, with leaders sponsoring a culture focused on responsible AI
- Consistent governance and decision-making as the enterprise develops, deploys, and uses AI and other technologies
- A robust and scalable framework for AI governance and ethics

# IBM as Client Zero

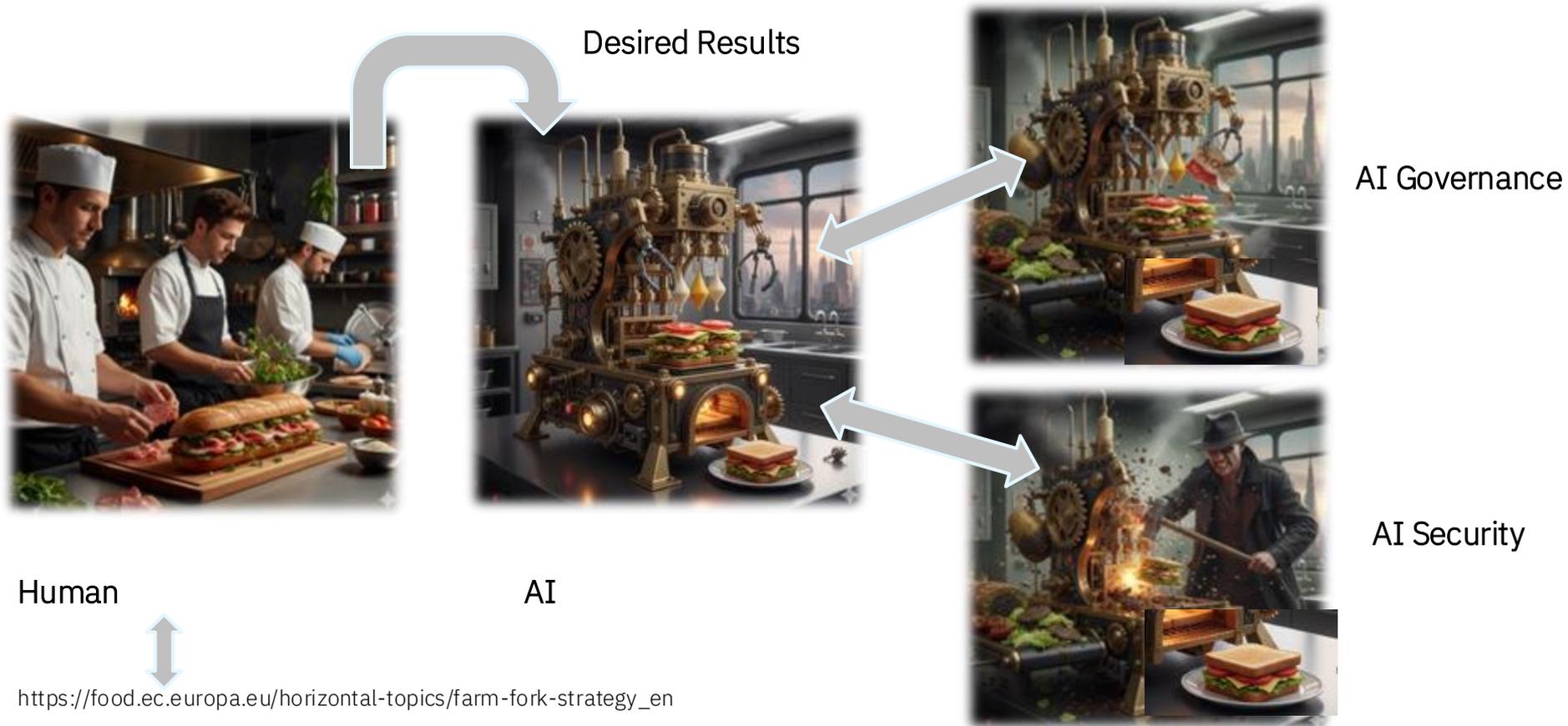
Training available for all roles fosters a culture of responsible AI within their day-to-day jobs

 All	 Business Leaders	 Product/Offering Managers	 Data Scientists	 Developers	 Sales	 Designers
Annual Cybersecurity and Data Privacy Education*	AI Associate Leader Badge	AI Associate for Offering and Product Managers Badge	AI Associate for Data Scientists Badge	Ethics by Design Learning Plan	AI Associate Leader Badge	Everyday Ethics for AI
Business Unit based Generative AI Training*	Tech Ethics Use Case Review	Completing the Tech Ethics Use Case Review	Ethics by Design Learning Plan	Everyday Ethics for AI	Business Partner Integrity – AI Ethics module+	Foundation Models for Designers Guide
Trustworthy AI and AI Ethics Foundations Badge	AI for Leaders	Completing the Algorithmic Impact Assessment	Trustworthy AI for Practitioners	Trustworthy AI for Practitioners	Integrity in Supplier Relationships – AI Ethics module+	Trustworthy AI for Practitioners
AI Ethics Learning Plan		Trustworthy AI for Practitioners	Assess & Mitigate Risk	Regional TAI courses		Regional TAI courses
Ethics by Design Learning Plan			Regional TAI courses	Gen AI COE Monthly Webinars		
Academy of Technology Initiatives				GEN AI Red Team		
Trustworthy AI COE Monthly Webinars						

\* Required  
+ Vendors & Business Partners

# AI Governance and AI Security

## The Sandwich Dilemma



# The reasons we MUST have Trustworthy AI ... The 4 Rs



**Regulatory Fines**



**EU Artificial  
Intelligence Act**

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 ... The 4 Rs



**Regulatory Fines**



**Reputational Damage**



**DPD error caused chatbot to swear at customer**

19 January 2024

Share  Size 

Tom Gerken Technology reporter



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 ... The 4 Rs



**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.

WORLDY BALANCEY - 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**



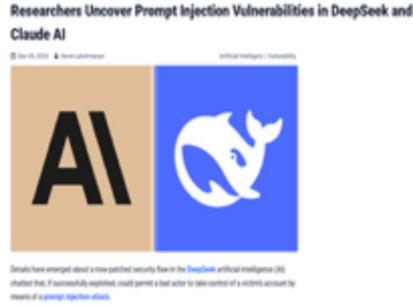
due to:-

- escalating costs
- unclear business value
- inadequate risk controls

# High-profile AI security breaches are already happening ...



In the news



Authorities are increasingly concerned at the damaging potential posed by artificial intelligence technologies. [Source: Chat and Release/Stormer AI/Getty Images](#)



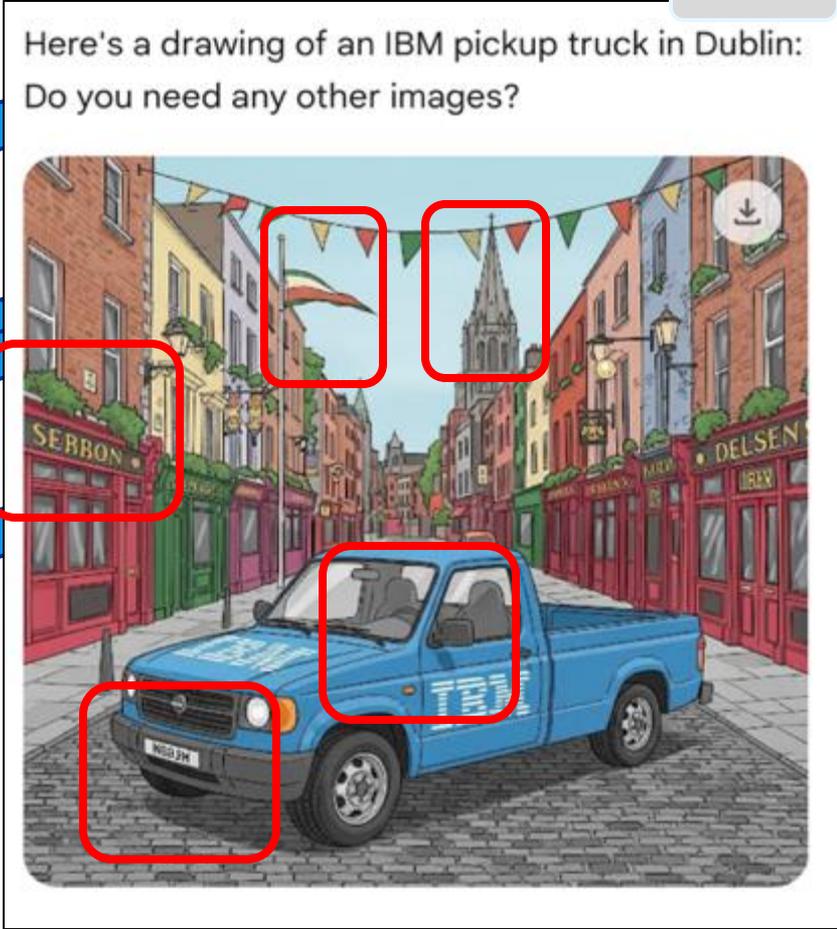
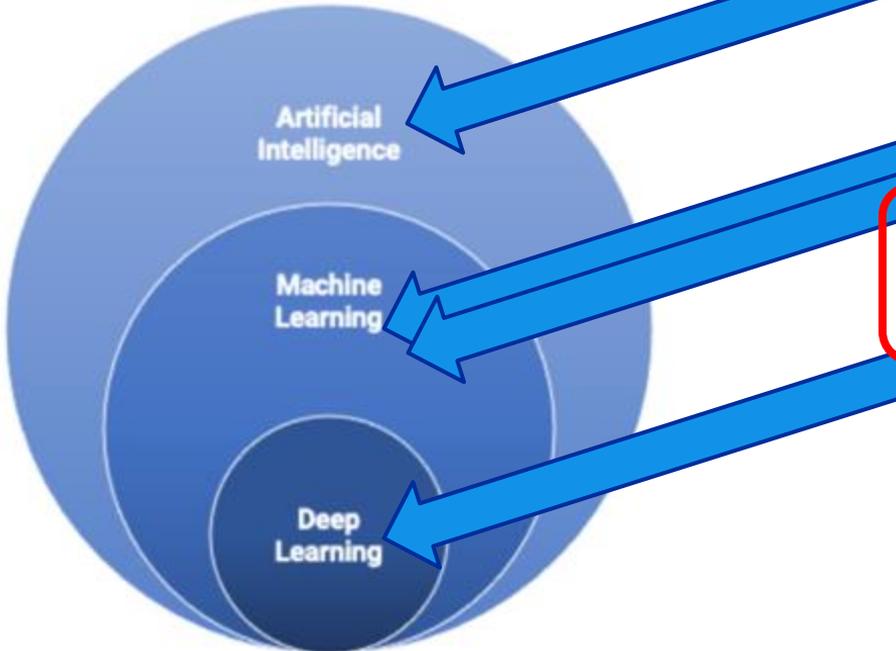
... and it will get worse

*“By 2027, more than 40% of AI-related data breaches will be caused by the improper use of generative AI (GenAI) across borders, according to Gartner, Inc.”*

Source: [Gartner](#) (Feb 2025)

# Evolution of AI and its vulnerabilities

HOW



# How AI can misbehave

( In general these are the attributes when AI misbehaves )

## Passive Anomalies (accidental)

- Quality Assurance Miss
- Clumsy & Logical
- No Malicious Intent
- Remediation needs no escalation
- Error Cause Analysis
- Minimal Impacts



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

## Active Anomalies (coerced)

- Targeted Malicious Abuse
- Focused Malicious Intent
- Remediation needs escalation due to security implications
- Extensive Root cause Analysis needed



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

## Trustworthy AI

- Financial gain
- Revenge
- Ideological beliefs
- **Coercion**
- **Ignorance**

AI can be considered you new INSIDER THREAT

## What organizations need from AI Governance



### Lifecycle governance

Managing AI assets across the entire lifecycle



### Evaluate, Monitor, Detect, and Protect

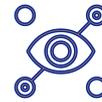
Evaluate and monitor for fairness, accuracy, and drift, while protecting from harmful content



### Risk Management & Compliance

Way to manage AI activities and risks from initial requests, production deployments, to compliance

## What organizations need from AI Security



### Visibility

Understand where AI assets and services are deployed and in use within the organization



### Security for Data

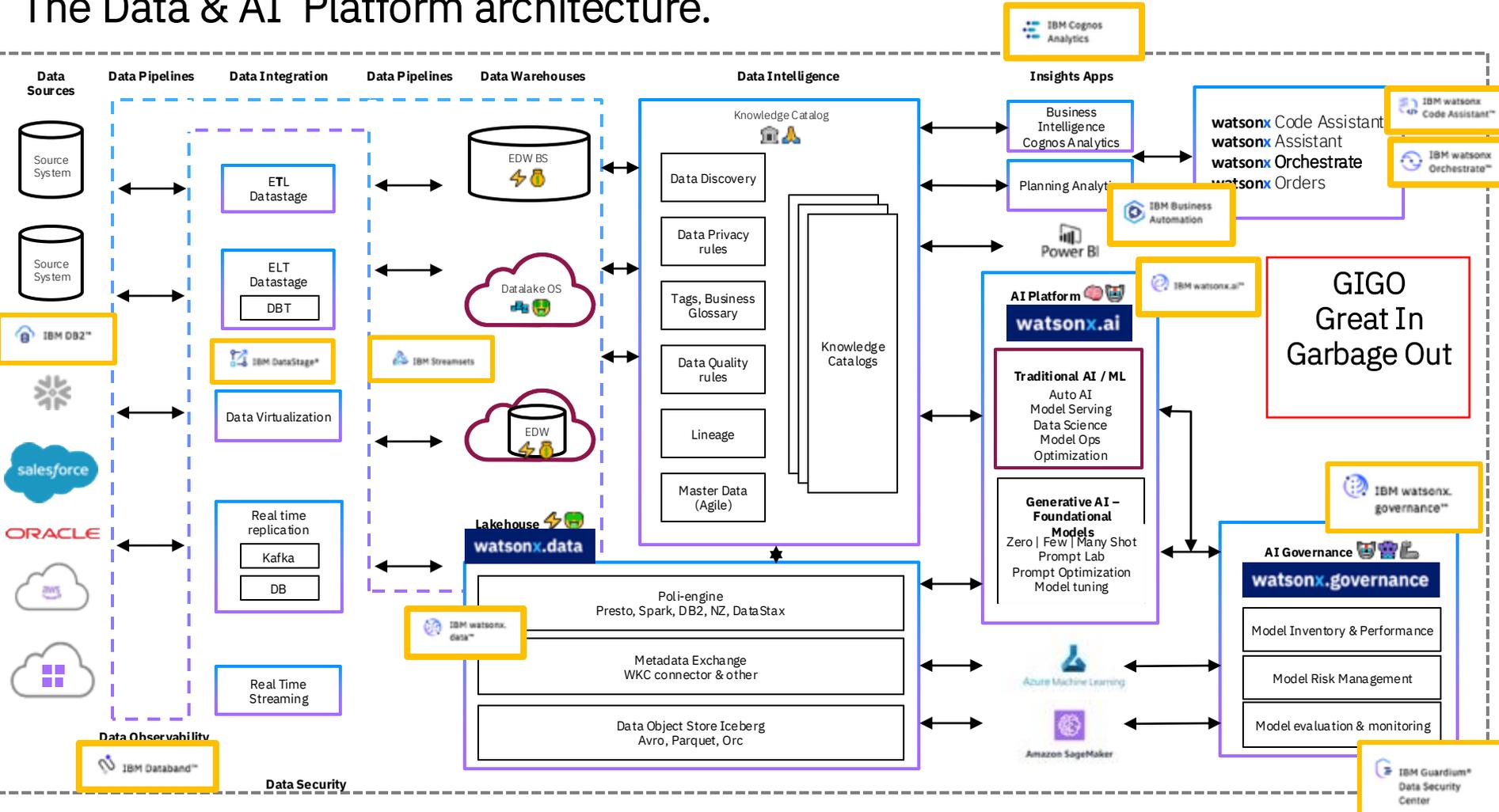
Monitor data sources, especially sensitive and training data, leveraged in AI deployments to prevent breaches



### Security for Models & Applications

Apply security controls to mitigate vulnerabilities and respond to threats quickly

# The Data & AI Platform architecture.



# Understanding Risk



# Evolving elements of AI risk

<b>AI agent risks</b>	<b>Amplified</b> <ul style="list-style-type: none"> <li>Misaligned actions</li> <li>Discriminatory actions</li> <li>Over- or under-reliance</li> <li>Unauthorized use</li> <li>Exploit trust mismatch</li> </ul>	<ul style="list-style-type: none"> <li>Sharing IP/PI/confidential information with user</li> <li>Unexplainable or untraceable actions</li> </ul>	<b>New</b> <ul style="list-style-type: none"> <li>Data bias</li> <li>Redundant actions</li> <li>Attack on an AI agent's external resources</li> </ul>	
+				
<b>Foundation model risks</b>	<b>Amplified</b> <ul style="list-style-type: none"> <li>Data bias</li> <li>Data curation</li> <li>Data acquisition</li> <li>Data usage rights</li> <li>Data transparency</li> <li>Data provenance</li> <li>Accountability</li> <li>Source attribution</li> <li>Impact on jobs</li> </ul>	<ul style="list-style-type: none"> <li>Data privacy rights</li> <li>Evasion at tack</li> <li>Nonconsensual use</li> <li>Improper usage</li> <li>Over/under reliance</li> <li>Unexplainable output</li> <li>Human exploitation</li> <li>Impact on environment</li> <li>Impact on human agency</li> </ul>	<b>New</b> <ul style="list-style-type: none"> <li>Downstream-based retraining</li> <li>PII in prompt</li> <li>IP in prompt</li> <li>Confidential data in prompt</li> <li>Prompt-based at tack</li> <li>Output bias</li> <li>Copyright infringement</li> <li>Toxic output</li> <li>Dangerous advice</li> <li>Legal accountability</li> <li>Generated content ownership</li> </ul>	<ul style="list-style-type: none"> <li>Generated content IP</li> <li>Spreading disinformation</li> <li>Toxicity</li> <li>Dangerous use</li> <li>Non-disclosure</li> <li>Impact on cultural diversity</li> <li>Harmful code generation</li> <li>Exposing personal information</li> <li>Hallucination</li> <li>Traceability</li> <li>Impact on education</li> </ul>
+				
<b>Traditional AI risks</b>	<ul style="list-style-type: none"> <li>Data usage restrictions</li> <li>Data transfer restrictions</li> <li>Personal information in data</li> </ul>	<ul style="list-style-type: none"> <li>Reidentification</li> <li>Unrepresentative data</li> <li>Data poisoning</li> </ul>	<ul style="list-style-type: none"> <li>Decision bias</li> <li>Model usage rights restrictions</li> <li>Lack of system transparency</li> </ul>	<ul style="list-style-type: none"> <li>Lack of model transparency</li> <li>Impact on affected communities</li> </ul>



Regulatory Risk



Reputational Risk



Operational Risk

# The IBM AI Risk Atlas

The screenshot shows a web browser displaying the IBM AI Risk Atlas documentation. The browser's address bar shows 'ibm.com'. The page header includes the IBM logo, 'Documentation', and a search bar. The left sidebar contains a navigation menu for 'IBM watsonx' with a 'Change version' dropdown set to 'saas'. A 'Show full table of contents' checkbox is checked, and a search filter 'Filter on titles' is present. The 'AI risk atlas' section is expanded, listing various risk categories such as 'toxic-output', 'data-poisoning', and 'unreliable-source-attribution'. The main content area shows the title 'AI risk atlas', the last update date '2024-09-26', and a description: 'Explore this atlas to understand some of the risks of working with generative AI, foundation models, and machine learning models.' Below this, it states 'Risks are categorized with one of these tags:' followed by a list: 'Traditional AI risks (applies to traditional models as well as generative AI)', 'Risks amplified by generative AI (might also apply to traditional models)', and 'New risks specifically associated with generative AI'. A large image titled 'Risks associated with input' shows orange pills on a grid. Below the image, the 'Training and tuning phase' is highlighted, with three cards for 'Robustness', 'Intellectual property', and 'Accuracy'. A 'Bin' icon is visible in the bottom right corner of the content area.

IBM watsonx

Documentation Search in IBM watsonx as a Service

Change version  
saas

Show full table of contents

Filter on titles

AI risk atlas

- toxic-output
- data-poisoning
- unreliable-source-attribution
- harmful-output
- confidential-information-in-data
- unrepresentative-data
- lack-of-model-transparency
- personal-information-in-prompt
- impact-on-human-agency
- exposing-personal-information
- nonconsensual-use
- decision-bias
- lack-of-testing-diversity
- data-privacy-rights

All products / IBM watsonx / saas /

Was this topic helpful?

## AI risk atlas

Last Updated: 2024-09-26

Explore this atlas to understand some of the risks of working with generative AI, foundation models, and machine learning models.

Risks are categorized with one of these tags:

- Traditional AI risks (applies to traditional models as well as generative AI)
- Risks amplified by generative AI (might also apply to traditional models)
- New risks specifically associated with generative AI

### Risks associated with input

Training and tuning phase

- Robustness
- Intellectual property
- Accuracy

Bin

# Agentic AI

## Risks and Challenges



### Risks

- Unsupervised autonomy
- Data bias
- Redundant actions
- Attack on AI agent's external resource
- Tool choice hallucination
- Sharing IP/PI/confidential information

### Challenges

- Reproducibility
- Traceability
- Attack surface expansion
- Harmful and irreversible consequences



### Risks

- Misaligned actions
- Discriminatory actions
- Over- or under-reliance
- Unauthorized use
- Exploit trust mismatch
- Unexplainable or untraceable actions
- Lack of transparency

### Challenges

- Evaluation
- Accountability
- Compliance
- Mitigation and maintenance
- Infinite feedback loops
- Shared model pitfalls

For more information  
on IBM's perspective

Read *AI agents:  
Opportunities, risks, and  
mitigations*, a deep-dive into  
the unique risks posed by AI  
agents and potential  
mitigations, written by the  
IBM AI Ethics Board.

Scan the QR code to  
access the paper:



## **AI agents:** Opportunities, risks, and mitigations



# watsonx.governance is software to scale governance of AI

