

Intelligent Maintenance Conference (IMC) 2024
Lausanne, Switzerland



Responsible AI Development and Legal Compliance: *Navigating the New Landscape of AI Regulation*

Rialda Spahic, PhD
Task Manager of Responsible AI, Equinor

NEWS ARTICLE | 1 August 2024 | Directorate General for Communication | 2 min read

AI Act enters into force



On 1 August 2024, the European Artificial Intelligence Act (AI Act) enters into force. The Act aims to foster responsible artificial intelligence development and deployment in the EU.



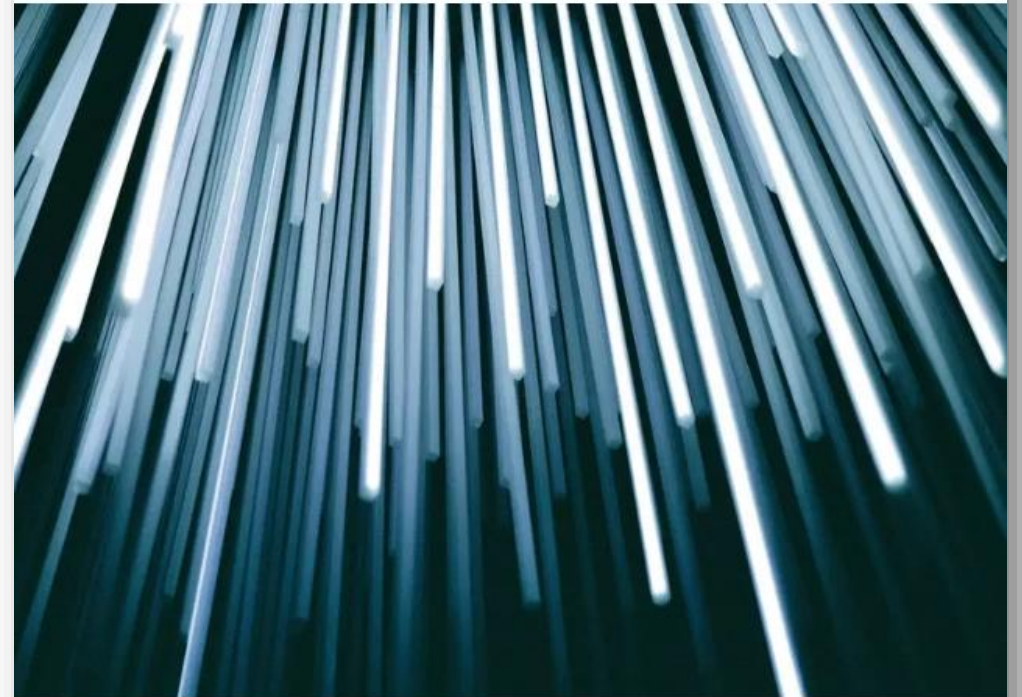
WORLD
ECONOMIC
FORUM

Join us

Sign in

Why we need to care about responsible AI in the age of the algorithm

Mar 22, 2023



Businesses need to take responsible AI seriously to remain competitive and avoid liability.

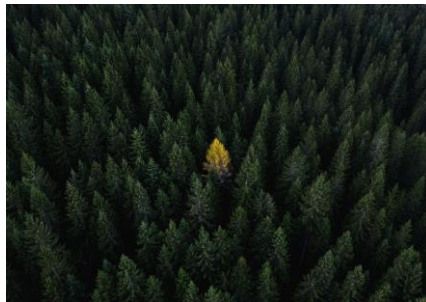
Image: Unsplash/Christopher Burns

Ayesha Gulley

Public Policy and Governance Associate,
Holistic AI



Uncertainty and risk mitigation



Environmental impact



AI Race



Security risks



Bias and fairness



Change in workplace

Beyond legal compliance

Why Responsible AI?

ML CO₂ Impact

Menu

ML CO₂ IMPACT

Machine Learning has a carbon footprint.

We've made a tool to help you estimate yours:

- 1 Compute your GPU's carbon emissions
- 2 Push for more transparency in our field by including the results in your publication (research paper, blog post etc.)
- 3 Install **codecarbon** to integrate carbon estimations in your Python workflow.

<https://mlco2.github.io/impact/>

AI Incident Database

Search over 3000 reports of AI harms

Search

Discover

IncidentDatabase.AI

Incident 701: American Asylum Seeker John Mark Dougan in Russia Reportedly Spreads Disinformation via AI Tools and Fake News Network

"Once a Sheriff's Deputy in Florida, Now a Source of Disinformation From Russia"

nytimes.com 2024-05-29

A dozen years ago, John Mark Dougan, a former deputy sheriff in Palm Beach County, Fla., sent voters an email posing as a county commissioner, urging them to oppose the re-election of the county's sheriff. He later masqueraded online as a R...

Read More →

<https://incidentdatabase.ai/>

AI Risk Repository

A Comprehensive Database of Risks from AI Systems

What are the risks from Artificial Intelligence?

A comprehensive living database of over 700 AI risks categorized by their cause and risk domain.

Read preprint

Explore database

<https://airisk.mit.edu/>



equinor

Responsible AI Principles

for designing, implementing and using AI



Governance and Accountability



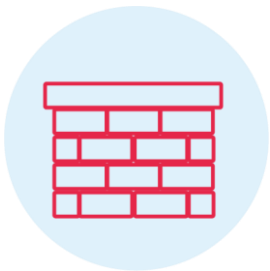
Safety and Security



Fairness and Openness



Human Oversight



Robustness and Resilience



Sustainability

AI Act enters into force



On 1 August 2024, the **European Artificial Intelligence Act (AI Act) enters into force**. The Act aims to **foster responsible artificial intelligence development and deployment in the EU**.



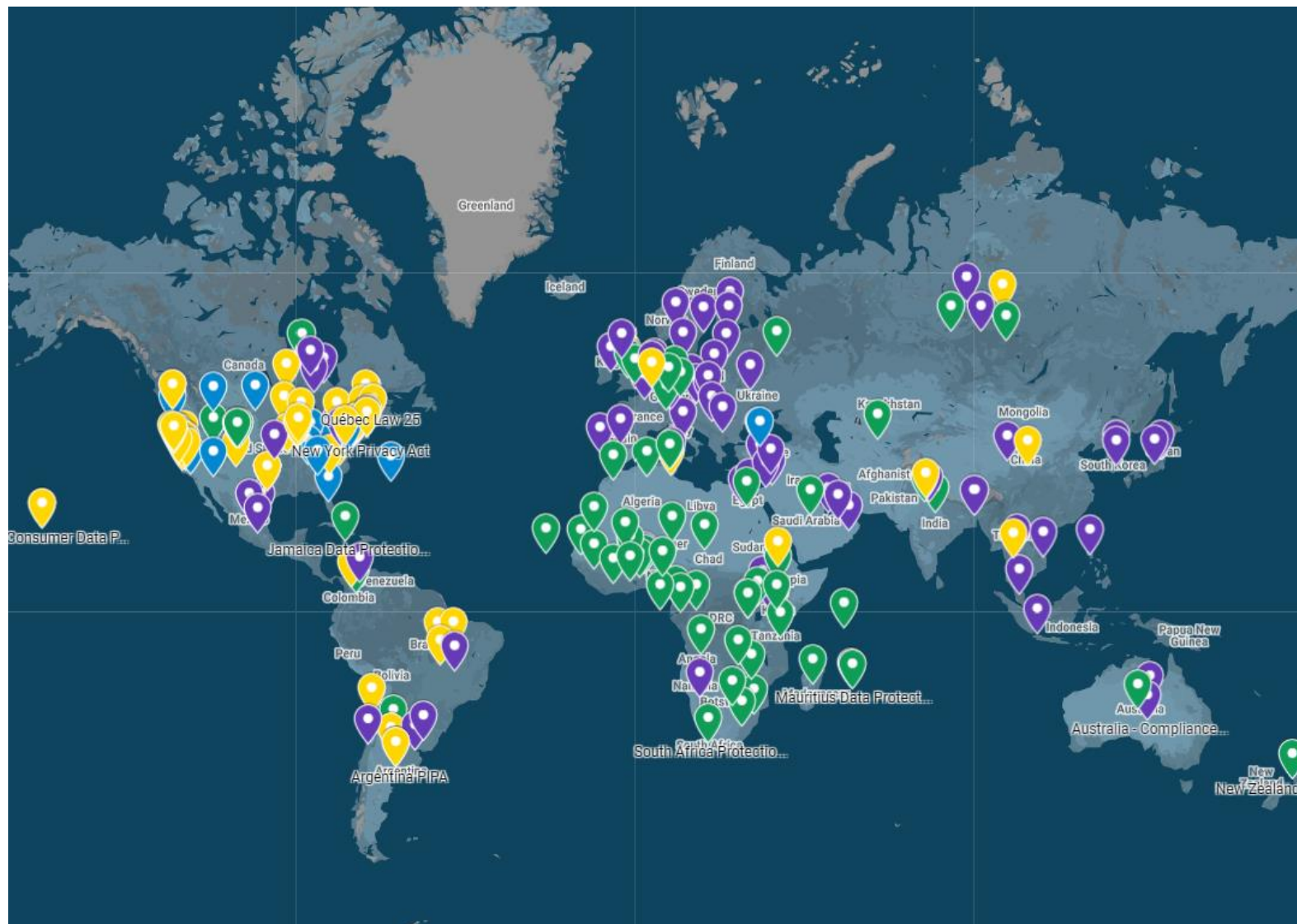


How to navigate this new landscape?

Our legal
obligations

Our obligations as
responsible members of
society

A Snapshot of Regulatory Landscape



Map of Global AI Regulations (April 2024)

Source: Fairly AI Map of Global Regulations April 23, 2024 [<https://www.fairly.ai/blog/map-of-global-ai-regulations>]



EU AI Act

GDPR, Data Protection Law

Regulation on Cybersecurity

Machinery Directive

Competition Law

Intellectual Property Law

Consumer Protection

Liability Law

Labor Law

Industry-specific regulations

Other laws, policies, regulations



AI Definition

EU AI Act, OECD



equinor

- ✓ ‘AI system’ means a machine-based system that is designed to **operate with varying levels of autonomy** and that **may exhibit adaptiveness after deployment**, and that, *for explicit or implicit objectives*, **infers**, from the input it receives, **how to generate outputs such as predictions, content, recommendations, or decisions** that **can influence** physical or virtual environments;

More
regulation

Unacceptable risk AI systems

e.g., Social scoring, subliminal manipulation

High-risk AI systems

e.g., Critical infrastructure, Employee management, Legal interpretation

Limited risk AI systems

e.g., AI Content generators, Chatbots

Low risk AI systems

e.g., Email spam filters

Less
regulation

More
regulation

Unacceptable risk AI systems

e.g., Social scoring, subliminal manipulation

High-risk AI systems

e.g., Critical infrastructure, Employee management, Legal interpretation

Limited risk AI systems

e.g., AI Content generators, Chatbots

Low risk AI systems

e.g., Email spam filters

Less
regulation

High-risk AI Systems

AI systems that can result in significant harm to people’s health, safety, fundamental rights or the environment.



Safety components for machinery



Safety or protective equipment used in potentially explosive environment



Lifts or safety components of lifts

Medical devices



Recreational crafts

Toys

AI systems used in the following areas are not automatically high-risk, unless they pose a significant risk:

Management and operation of critical infrastructure



Biometric identification and categorization of natural persons

Employment, worker management



Essential private and public services and benefits

Assistance in legal interpretation and application of the law



Education and vocational training

Law enforcement



Migration, asylum and border control management



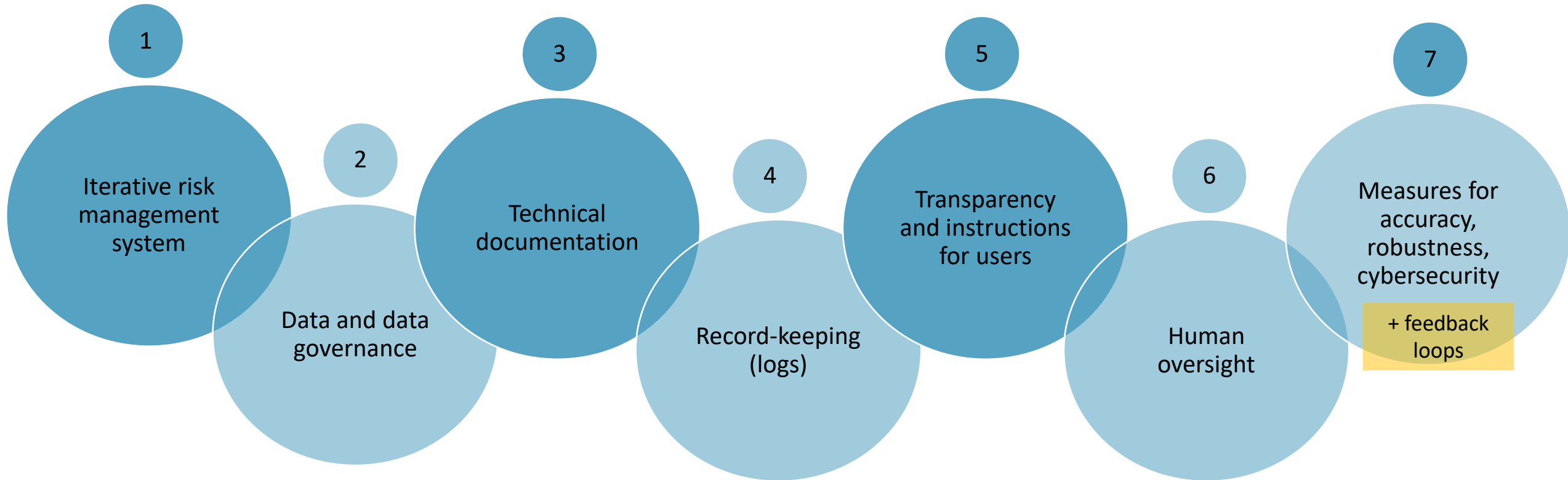
equinor

High-Risk AI System Criteria

Criteria used to assess whether an AI system poses a risk of adverse impact on fundamental rights:

- Intended **purpose** of AI system
- Potential extent of the **harm**
- The extent to which harmed persons are in a **vulnerable** position
- The extent in which the outcome of the system is **reversible**
- The extent in which **existing legislation** provides for effective **measures** to address and minimize the risks

Design requirements for High-risk AI Systems



EU AI Act



More
regulation

Unacceptable risk AI systems

e.g., Social scoring, subliminal manipulation

High-risk AI systems

e.g., Critical infrastructure, Employee management, Legal interpretation

Limited risk AI systems

e.g., AI Content generators, Chatbots

Low risk AI systems

e.g., Email spam filters

Less
regulation

EU AI Act

Non-compliance

Non-compliance fines up to 35 million EUR
or **7%** of a company's annual turnover

Preparing for compliance

- ✓ **AI Register/Database**
- ☐ AI Risk Management
- ☐ AI Governance Framework and Guidelines



AI Register



High-Risk AI Systems Database

Internal database listing AI systems
(In-house, Third Party, Vendor)
and their documentation.

EU AI Act:
Requirement to register all high-risk AI
systems in EU AI Database of High-Risk AI
Systems, along with documentation

Preparing for compliance

- ✓ AI Register/Database
- ✓ **AI Risk Management**
- ☐ AI Governance Framework and Guidelines



AI Risk Library, Risk Taxonomy

- First step towards risk assessments
- Internal AI risk landscape and a focus area heatmap
- Assist audit, procurement and product stakeholders

AI Incidents

Any outcome of the system that could cause harm



Attacks

Backdoors
Data Poisoning
Model extraction
...

Intentional abuse

Ethnic Profiling
AI-enhanced Cybercrime
Misuse of autonomous bots
...

Failures

Data Drift
Discrimination
Opacity
...

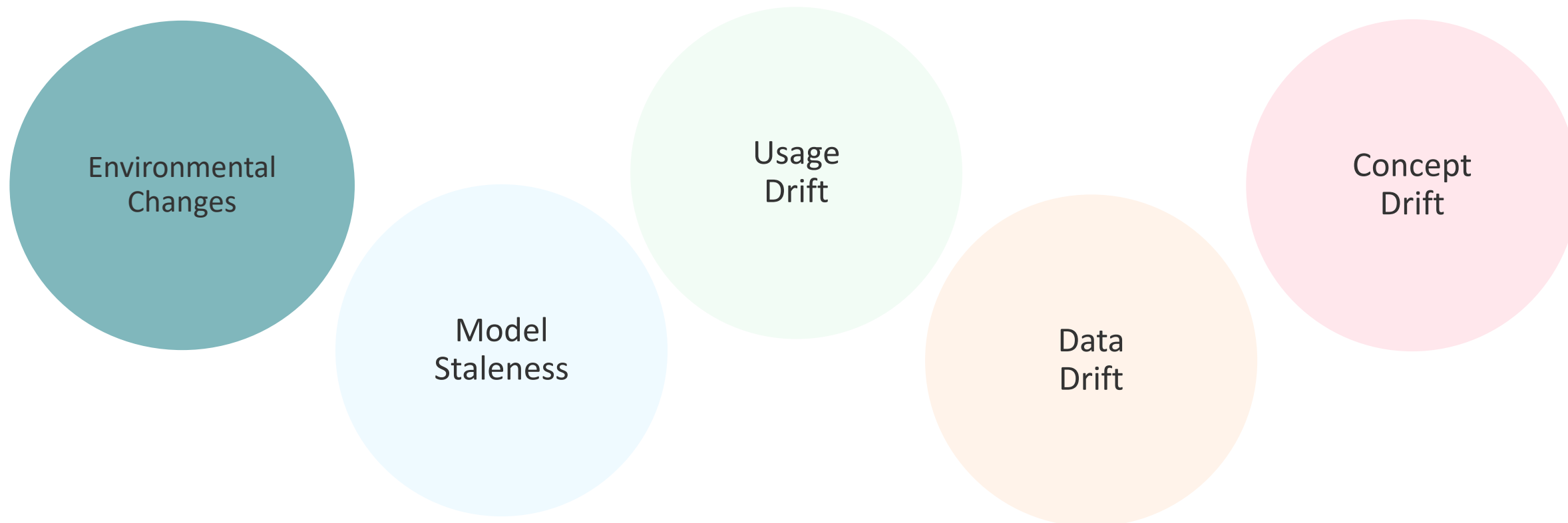
Example

AI Model Decay and Ripple Effects

An AI system in **development** may **not** be the same once *operational*.



Degradation of model's performance over time
Gradual or sudden



Model Staleness

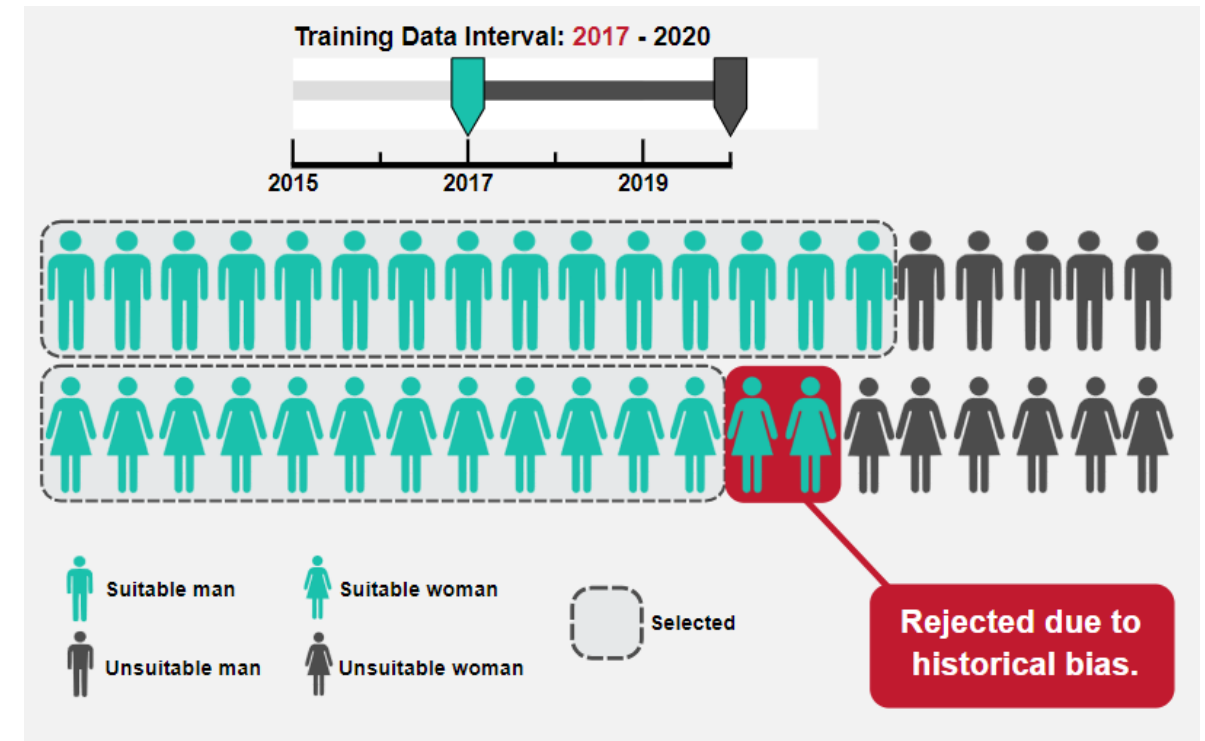
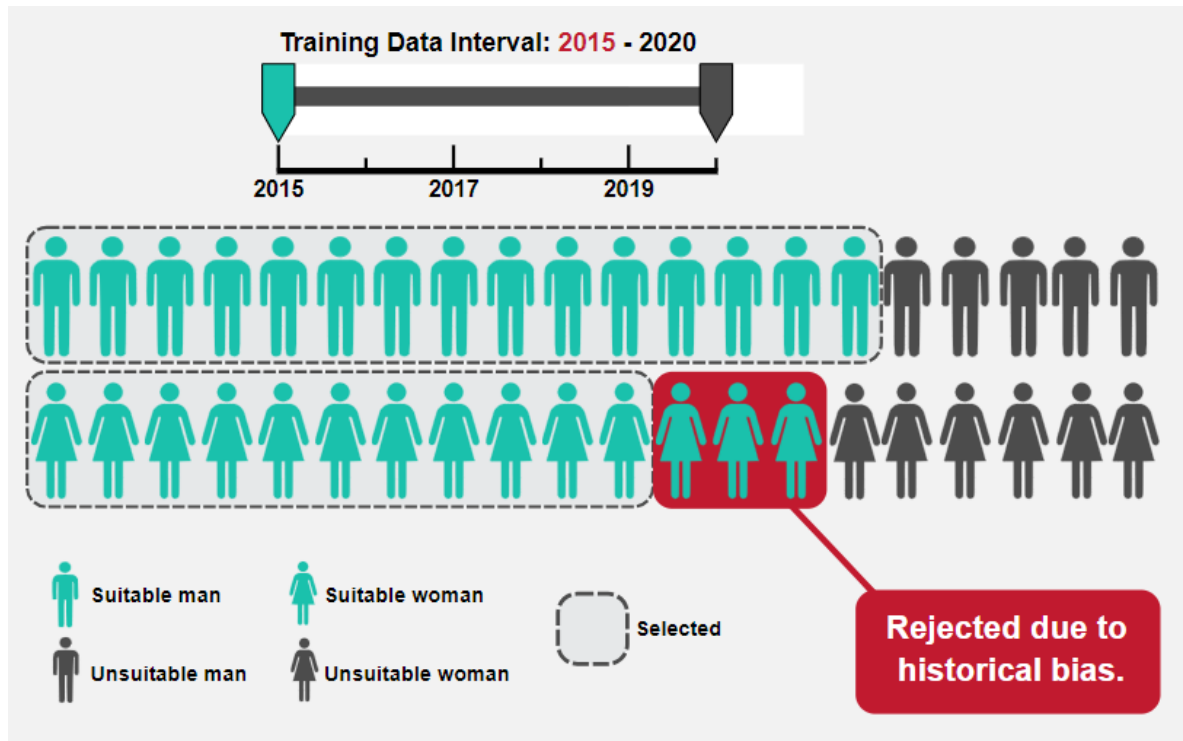
Example

AI Model Decay and Ripple Effects

Aging algorithms and outdated assumptions;
Data used to train models does not represent current reality;



equinor



Historical bias in AI systems

Figure source" Australian Human Rights Commission

<https://humanrights.gov.au/about/news/media-releases/historical-bias-ai-systems>, Accessed on 08.07.2024

Example

AI Model Decay and Ripple Effects



kDimensions

Data Drift

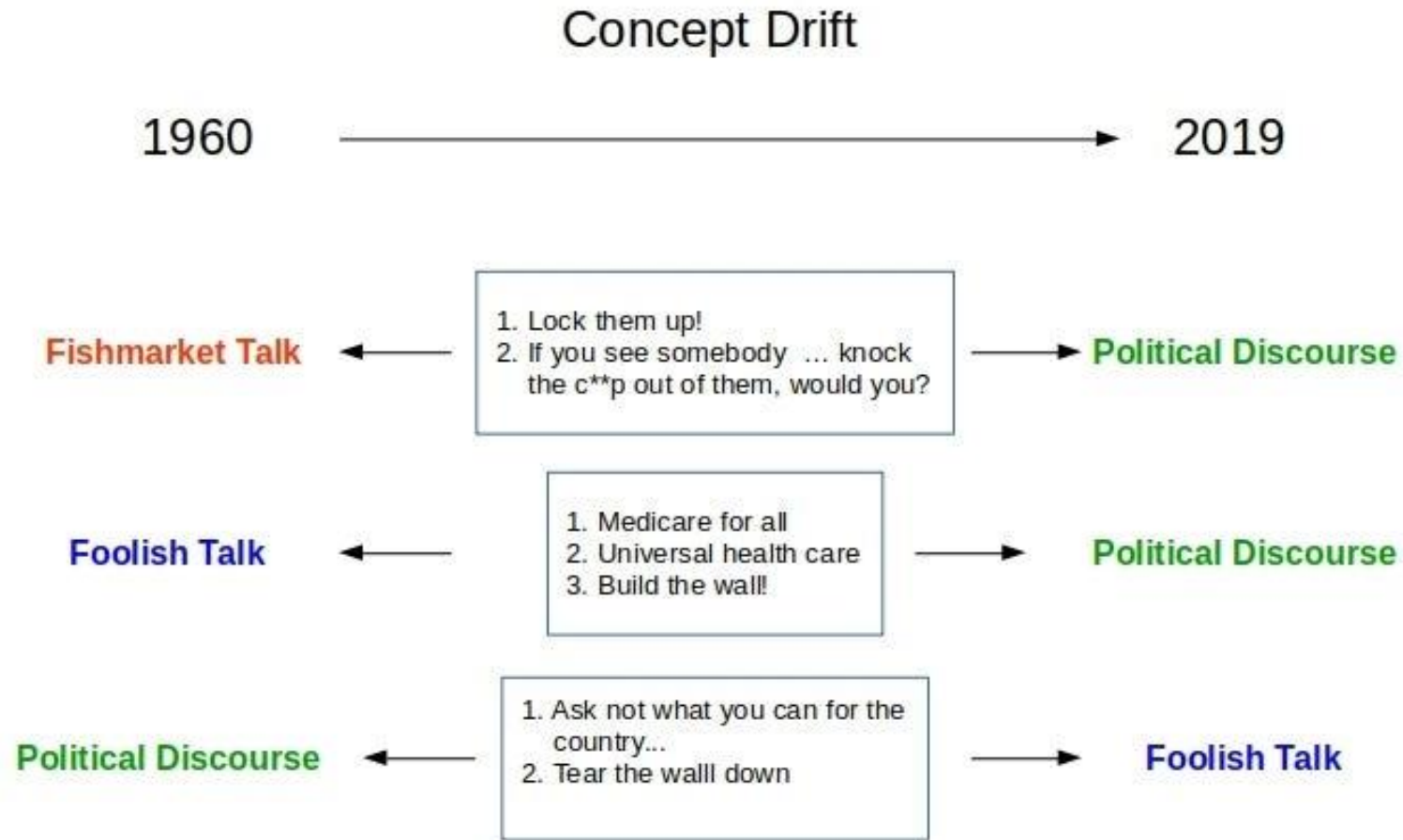
Change in statistical characteristics of the **input data**;
The input data has changed rendering the trained model irrelevant on new data;

Concept Drift

Changes in the underlying **relationships** in data;
Objective change = what we are trying to predict has changed;

Example

AI Model Decay and Ripple Effects



Data has not changed, but our interpretation and assigned class have

Mitigating risks

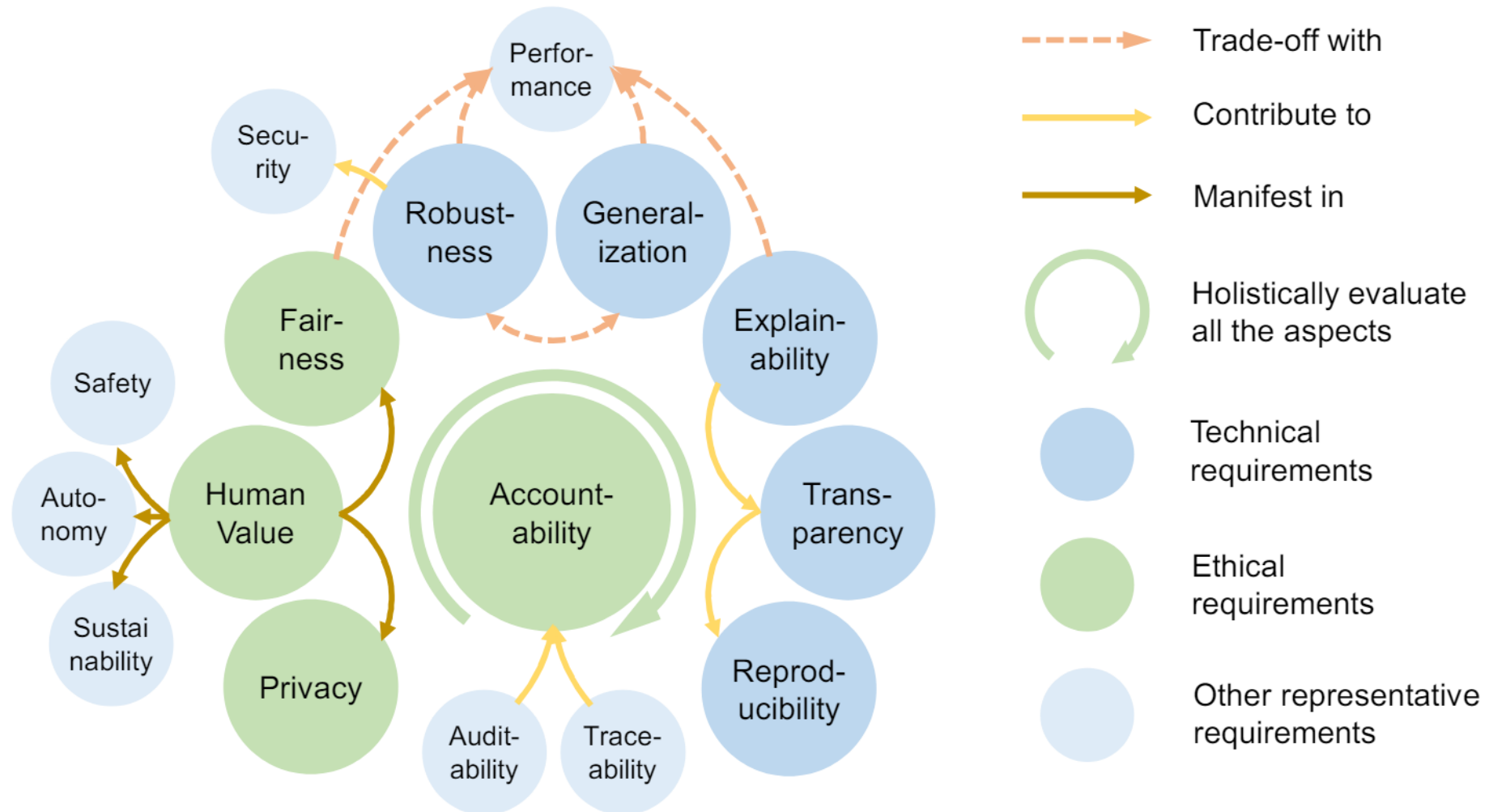
- EU AI Act: Risk Management System
- National Institute of Standards and Technology (NIST): AI Risk Management System
- ISO/IEC 23894:2023 Artificial intelligence - Guidance on risk management
- ISO/IEC 42001 Artificial intelligence – Management system
- US Federal Reserve SR 11-7 Model Risk Management
- ...

Risk analysis	Risk evaluation	Risk mitigation	Residual risk acceptability evaluation	Risk management report	Usage instructions, documentation, EU Register
<p>Gather data:</p> <ol style="list-style-type: none">1. Initial project plan2. User requirements3. Data from previous projects / versions4. Data on similar systems5. Standards	<p>Estimate risk levels for every hazardous scenario and identify scenarios with high risk levels</p>	<p>Identify potential mitigation measures including changes to user interface as well as other changes</p>	<p>Assess and document residual risk evaluation and justify why these are acceptable</p>	<p>Provide a risk management report</p>	<p>Risk management report as a part of Technical documentation</p>
<p>Draft statement:</p> <ol style="list-style-type: none">1. Intended users2. Expected levels of user experience3. Use environment4. System functionalities5. Anticipated effect on the environment6. Operating principles	<p>Do a thorough analysis of high-risk scenarios, identify correlations and causations to help design adequate mitigation measures</p>	<p>Identify necessary additional safety information in the 'Instructions for use'</p>	<p>Consider conducting AI sandbox testing</p>	<p>Design monitoring processes and define thresholds for re-assessment</p>	<p>Draft appropriate 'Instructions of use' based on risk assessment</p>
<p>Examples of use scenarios</p>	<p>Do a cross-impact analysis of all identified risks</p>	<p>Assess effectiveness of individual measures</p>	<p>Consider conducting real-world testing</p>	<p>Define fixed timeframes for mandatory re-assessments</p>	<p>Design training courses / materials if necessary to mitigate risks</p>
<p>Document identified:</p> <ol style="list-style-type: none">1. Foreseeable use errors2. Foreseeable intentional misuse and abnormal use	<p>Re-evaluate all high-risk scenarios</p>	<p>Document residual risks and corresponding risk levels</p>		<p>Update report as necessary</p>	<p>Submit risk management report when registering the system in EU Register</p>
		<p>Formative evaluation of residual risks</p>			
		<p>Summative evaluation of residual risks</p>			

Moving beyond obvious risk

AI Trade-Offs

Prioritization of Responsible AI Practices



Cultural competencies for AI Risk Management



Organizational accountability

- written policies and procedures
- effective challenge
- accountable leadership

Organizational processes

- forecasting failure modes and knowing past failures
- deliberating on **who** (customers, stakeholders),
what (well-being, opportunities),
when (frequently, over long period of time),
how (immediate response, altering processes)

Culture of effective change

- diverse and experienced teams
- challenging, evaluating, assessing at each step of the AI lifecycle

Preparing for compliance

- ✓ AI Register/Database
- ✓ AI Risk Management
- ✓ **AI Governance Framework and Guidelines**



Responsible AI guidelines and best practices

Guidelines for each major process in AI lifecycle aligned with regulatory obligations, internal governance and best practices

- Ideation
- Development
- Procurement
- Deployment
- Operations and monitoring
- Use



Lessons learned

AI is used to do things faster, more consistent
and at scale



AI incidents can occur **fast, consistently, and at scale**



- ☐ Performance metrics
- ☐ Monitoring and continuous monitoring
- ☐ Regular model retraining
- ☐ Adaptive algorithms

Governance → Drastic increase in materiality of
AI systems must be
followed by drastic increase of AI governance
and safety measures



- ☐ Anticipate changes
- ☐ Feedback loops
- ☐ Human in the loop
- ☐ Continuous upskilling

Way ahead and what to expect



EU AI Act Timeline

- 6 months for prohibited AI systems.
- 12 months for GPAI.
- 24 months for high-risk AI systems under Annex III.
- 36 months for high-risk AI systems under Annex II.
- EU AI Office Codes of practice shall be ready 9 months after entry into force

EU AI Office work in progress

- facilitate compliance process
(for high-risk AI systems- no later than 18 months after entry into force)
- facilitate the drawing up of codes of practice
- facilitate detection and labelling of artificially generated or manipulated content
- The EU Commission will provide guidelines and examples of high-risk and non-high-risk AI systems and add or remove conditions for high-risk classification based on evidence.

Intelligent Maintenance Conference (IMC) 2024
Lausanne, Switzerland



Thank you

**Responsible AI Development and Legal Compliance:
*Navigating the New Landscape of AI Regulation***

Rialda Spahic, PhD
Task Manager of Responsible AI, Equinor

Image source: Equinor Media Bank