# Trustworthy AI in action
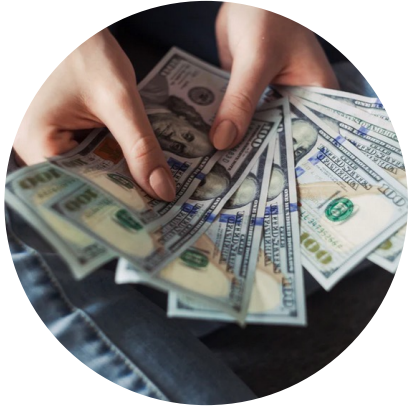
July 26, 2022

John Thomas
VP & DE, IBM Expert Labs
@johnjaithomas

Ashley Casovan
Executive Director
Responsible AI Institute
@ResponsibleAI

IBM

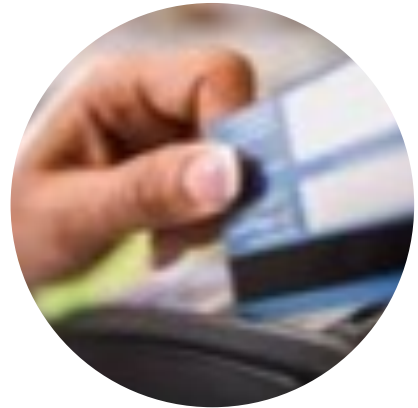# With AI increasingly powering critical workflows...



credit



employment



customer management
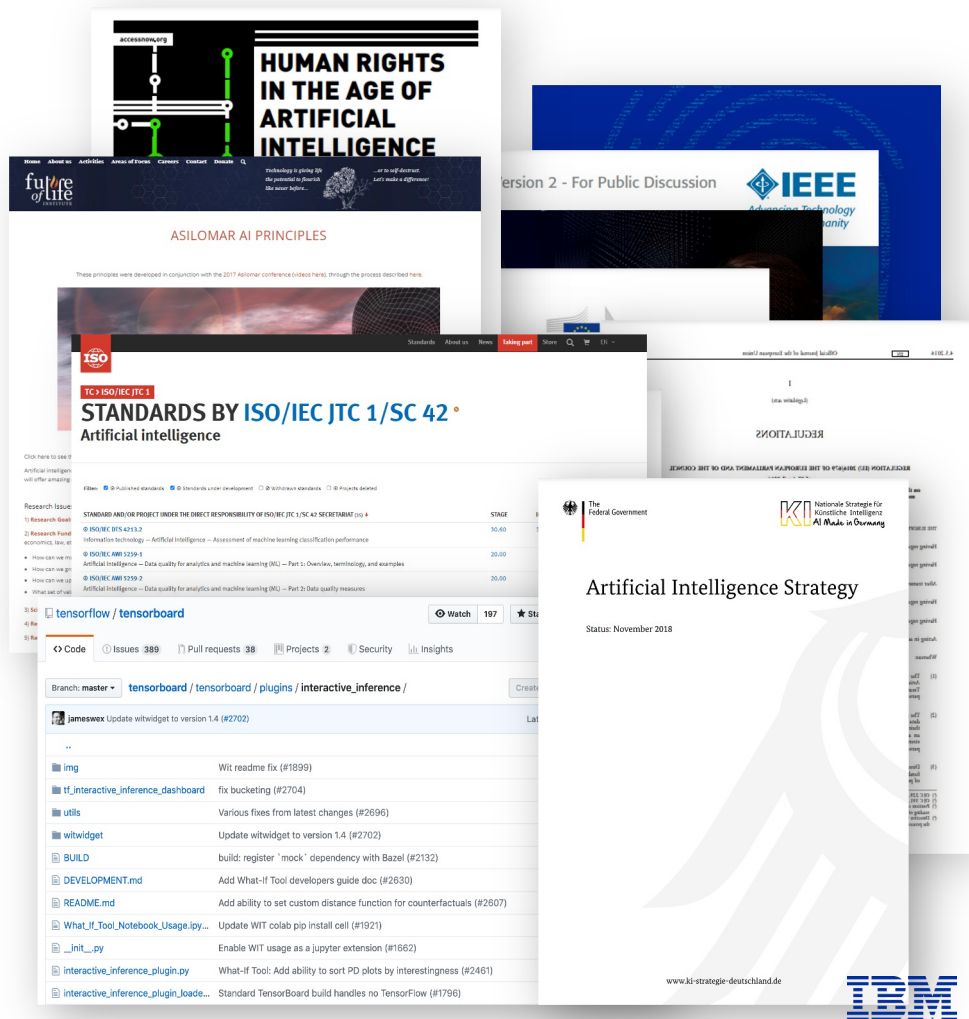


fraud

*...**trust** is essential*

# Simplifying how to implement responsible AI

"Only about a quarter (28%) of citizens are willing to trust AI systems in general. Two out of five citizens are unwilling to share their information or data with an AI system and a third are unwilling to trust the output of AI systems."

- University of Queensland and KPMG, 2021

"Fewer than 20% of executives strongly agree that their organizations' practices and actions on AI ethics match (or exceed) their stated principles and values."

- IBM and Oxford Economics, 2021

# Organizations must consider Regulatory Compliance



**Sarbanes-Oxley Act**

## USA

2021 – National AI advisory committee

2022—Algorithmic Accountability Act of 2022

2022 – American Data Privacy and Protection Act

## Canada

2017- National AI Strategy

2020—Directive on Automated Decision Making

2022—Artificial Intelligence and Data Act

## European Union

2018— Coordinated Plan on AI

2021 – Draft AI Act
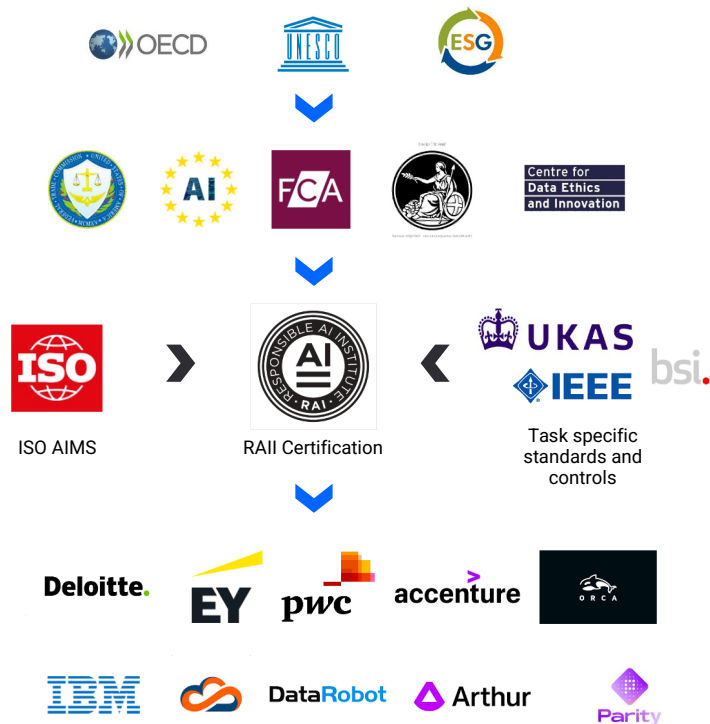
## United Kingdom

2021— CDEI AI Assurance Guide

## Standards Bodies

National standards and accreditation bodies are working on AI specific standards and frameworks to help implement responsible AI (eg. ISO, NIST, SCC, CEN-CENELEC, BSI, UKAS, ANSI, CETA, IEEE, etc.)

# AI Regulatory and Standards Landscape

| AI Principles | <ul><li>International NGOs</li><li>Corporate values</li><li>ESG objectives</li></ul> |
| --- | --- |
| AI Regulations | <ul><li>National</li><li>State, Regional, and Local</li><li>Corporate policies</li></ul> |
| Standards, Certifications, and Industry Best Practices | <ul><li>Accreditations</li><li>Management Standards</li><li>Use case & function specific certifications</li><li>Task specific standards & controls</li><li>Industry best practices</li></ul> |
| Evaluations | <ul><li>Point-in-time audits<ul><li>Manual</li><li>Semi-automated</li></ul></li><li>Ongoing monitoring<ul><li>Statistical evaluations</li><li>Data quality evaluations</li><li>Automated policy evaluations</li><li>Document tool chain</li></ul></li></ul> |

ISO AIMS

RAII Certification

Task specific standards and controls

# We need a multidisciplinary, multidimensional approach to trustworthy AI
## From principles to actions



**what should be done**
principles, values, norms, laws, regulations



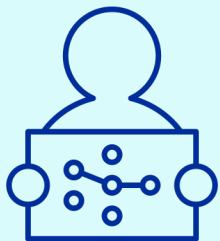**how to instrument it**
techniques, algorithms, software, best practices



**how to operationalize it**
mechanisms, systems, and processes to keep AI trustworthy

IBM

# What does it take to Trust a decision made by an AI?
## *We started from these HUMAN-CENTRIC questions*

| | | | | |
|---|---|---|---|---|
| Is it easy to understand? | Is it fair? | Did anyone tamper with it? | Is it accountable? | Does it safeguard data? |
| **EXPLAINABILITY** | **FAIRNESS** | **ROBUSTNESS** | **TRANSPARENCY** | **PRIVACY** |
| Easy to understand outcomes/decisions | Impartial and addressing bias | Handle exceptional conditions effectively | Open to inspecting facts and details | High integrity data & business compliance |
| *Why did the AI arrive at an outcome? When would it have been different?* | *Are privileged groups at a systematic advantage compared to other groups?* | *Can we evaluate and defend against a variety of threats?* | *Can we increase understanding of why and how AI was created?* | *How do we ensure owners retain control of data and insights?* |

IBM

# AI governance enables trustworthy AI

## Foundational for strategy and execution of AI solutions

| Strategy | Planning | Development and deployment | Operate | Monitor + portfolio management |
|---|---|---|---|---|
| **Who?** Business, AI ethics board, data/AI leaders, people responsible for internal policy and regulations (CPO etc.) | **Who?** Business, AI ethics board, data/AI leaders, people responsible for internal policy and regulations (CPO etc.), ecosystem | **Who?** Dev teams, IT leaders, CDO, people responsible for software and data scientists | **Who?** IT leaders, MLOps teams | **Business outcomes** Who? Business leaders + MLOps teams  **Governance of models** Who? Business leaders + MLOps teams |

**AI strategy**
Humans in the loop from the beginning

**AI governance tools and process**
Ensuring governance of AI models through technology and data

Decide and drive the AI strategy for the organization. Establish AI policies for the organization. (This may include AI principles, regulations, laws, etc.) => AI Ethics Board

Enable data collection and transparent reporting to make needed information available to all stakeholders.

Encode the policies into business rules, guidelines and transparent reporting mechanism. Determine appropriate guardrails and parameters.

IBM

# RAII Implementation Framework Dimensions and Sub-Dimensions



**1** **Systems Operations**

1.1 System Scope and Function
1.2 Human-in-the-Loop
1.3 Model is Fit for Purpose
1.4 Data Relevance and Representativeness
1.5 Data Quality

**2** **Explainability & Interpretability**

2.1 Communication About the Outcome
2.2 Notification
2.3 Recourse
2.4 Understanding the AI System's Decisions or Functions

**3** **Accountabilty**

3.1 Organizational Governance
3.2 Team Governance

**4** **Consumer Protection**

4.1 Transparency to the User and Data Subject
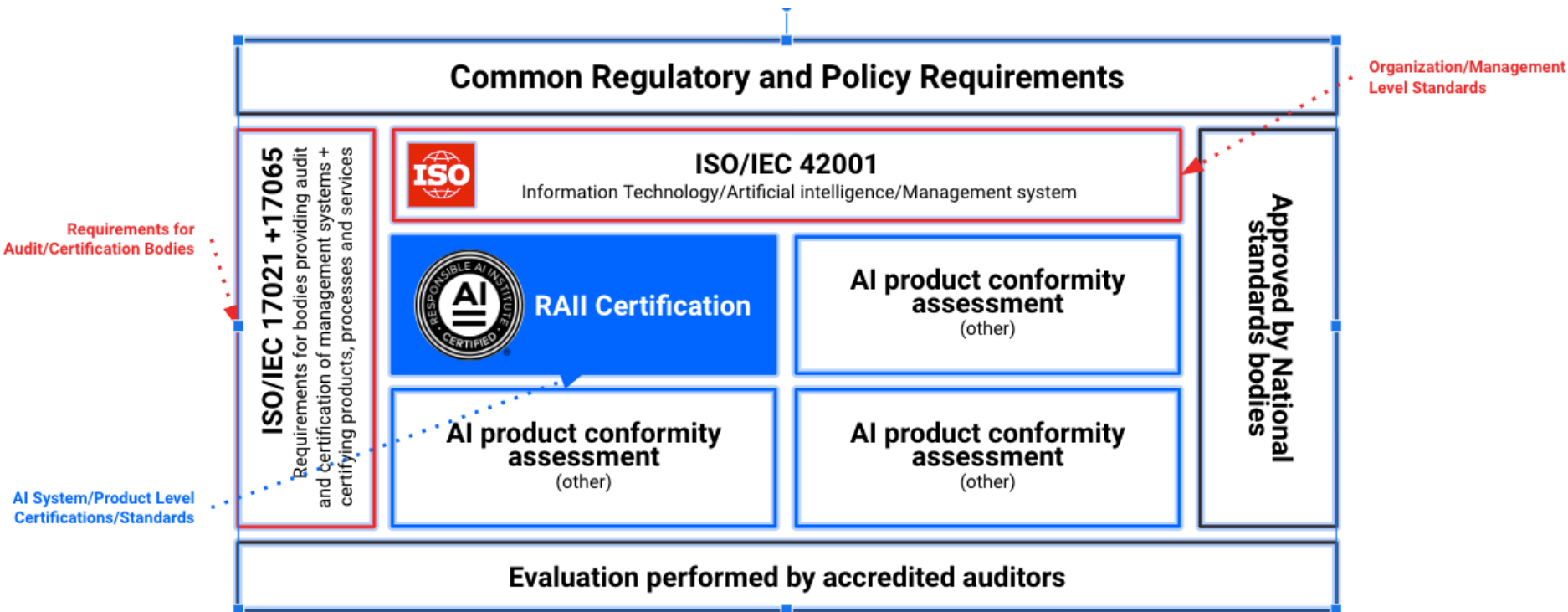4.2 Harms to Individuals
4.3 Protections

**5** **Bias & Fairness**

5.1 Bias Impacts
5.2 Bias Training
5.3 Bias Testing

**6** **Robustness**

6.1 Data Drift
6.2 System Acceptance Test is Performed
6.3 Contingency Planning

Responsible
Artificial Intelligence
Institute

# Complying with AI regulations

# Sample Control Development

## OECD AI Guiding Principle on Fairness (Global Level)
AI systems should be designed in a way that respects the rule of law, human rights, democratic values and diversity, and they should include appropriate safeguards – for example, enabling human intervention where necessary – to ensure a fair and just society.

## Enterprise AI Policy (Industry Level)
**2. Ensure Fairness in AI Systems**
2.1 Identify unwanted bias.   2.2. Test for unwanted bias.   2.3 Perform bias training.   2.4 Ensure recourse is implemented

## Client Bias Guideline (Company Level)

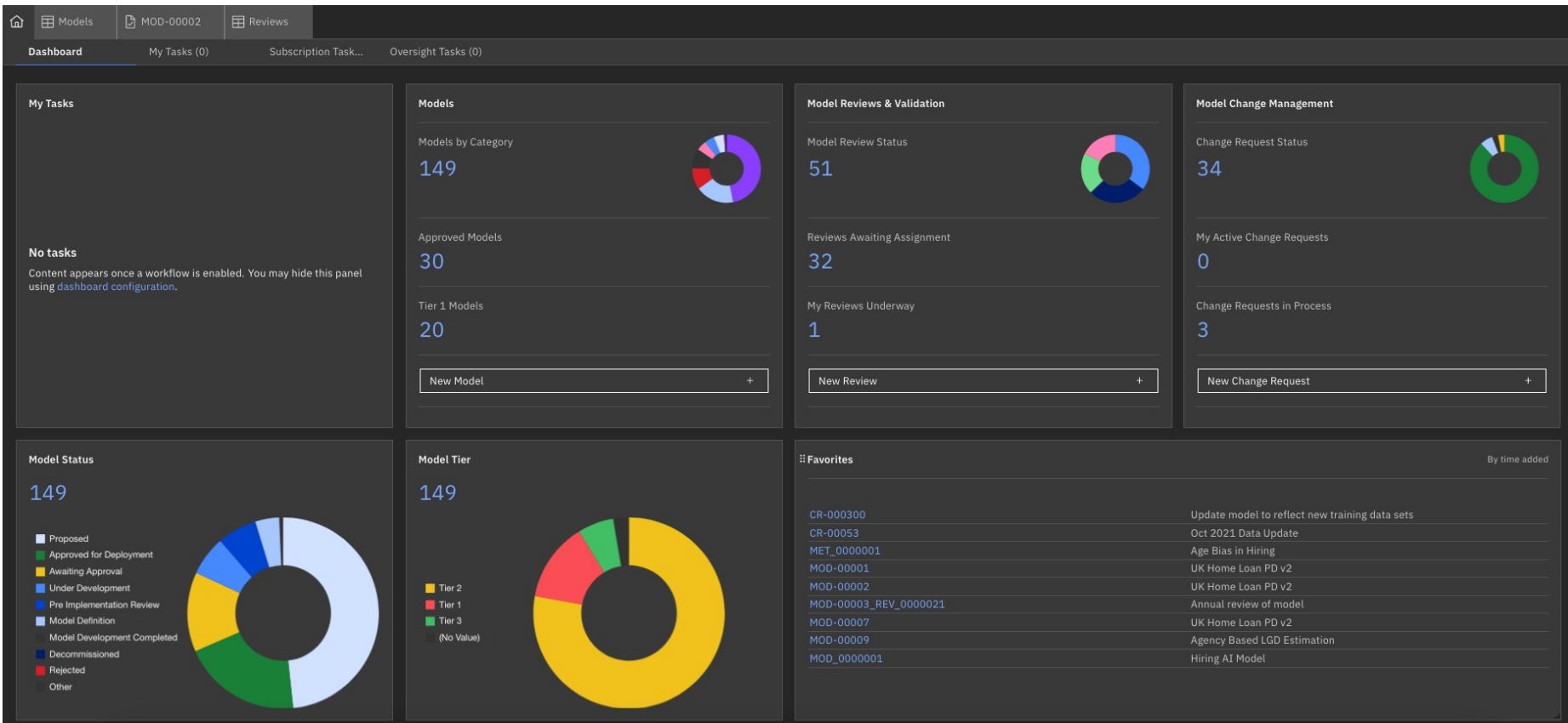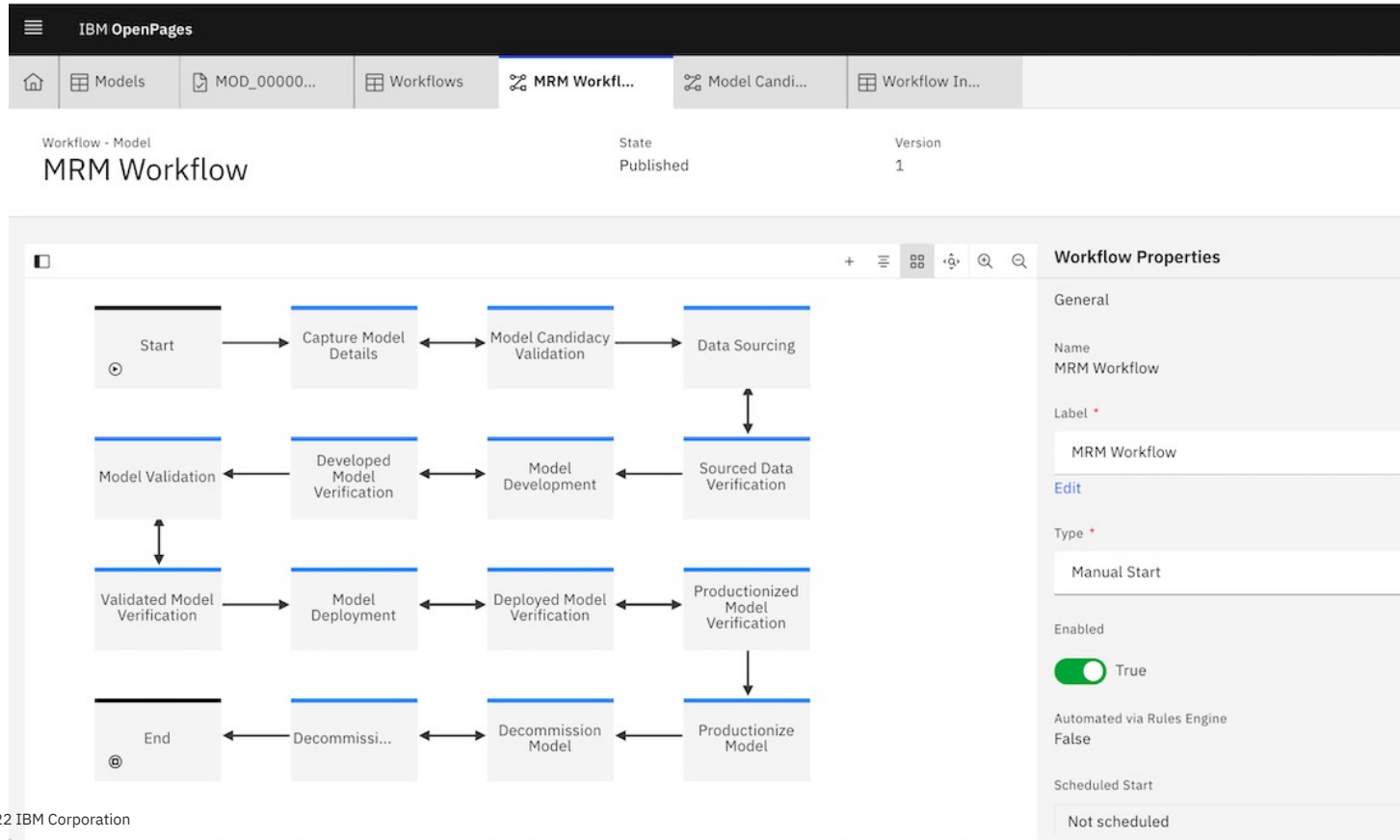| | | | |
|---|---|---|---|
| 2.1 Complete a harms mapping to understand the unintended consequences and identify mitigation measures | 2.2 Perform appropriate bias test for all protected classes, ensure that it meets the acceptable threshold. | 2.3 Require bias training for all employees involved in the production and deployment of an AI system | 2.4 All AI systems must have a recourse plan in the case of the system not working as intended |

## Bias Control (AI System/UC Level)

2.2 Perform demographic parity and equalized odds test for protected classes. Acceptable result is 0.85.

# Enterprise Governance View of AI/ML models



© 2022 IBM Corporation

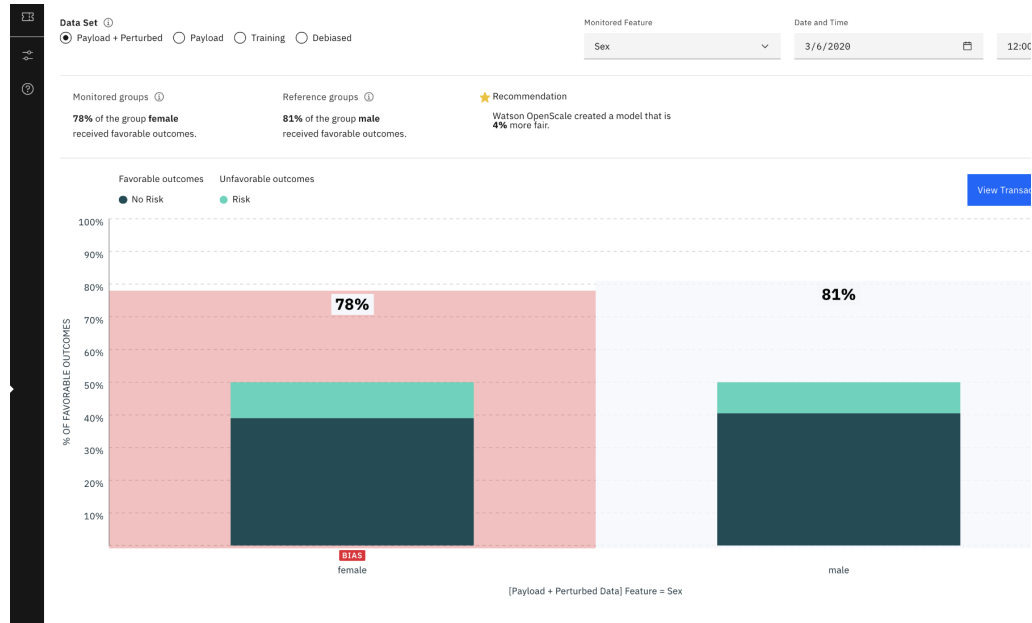# Enterprise workflows provide consistency with flexibility

# Example of a runtime guardrail
*Bias Monitoring and Mitigation with Watson OpenScale*



Setup ongoing monitoring of deployed model
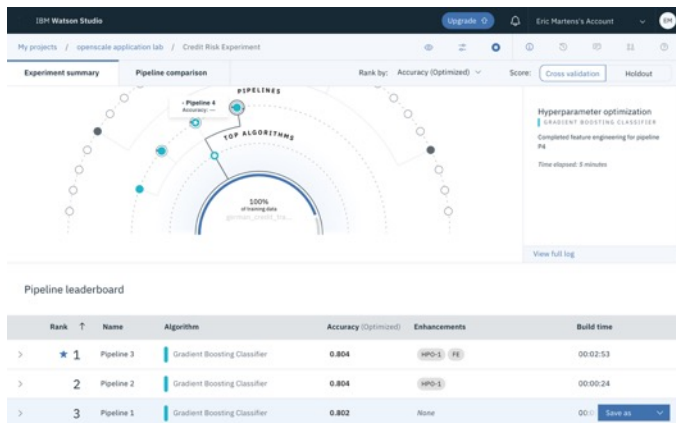Define monitored and reference groups
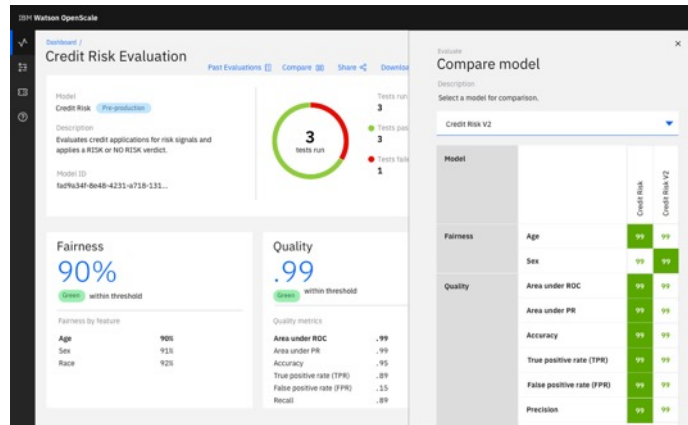
Calculate Disparate impact Value

78% of the monitored group (female) have a favorable output
81% of the reference group (male) get a favorable output

Disparate impact Value: 96%
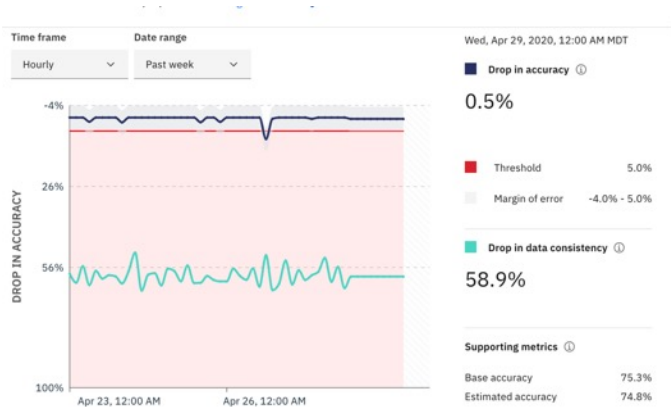Mitigation based on policy
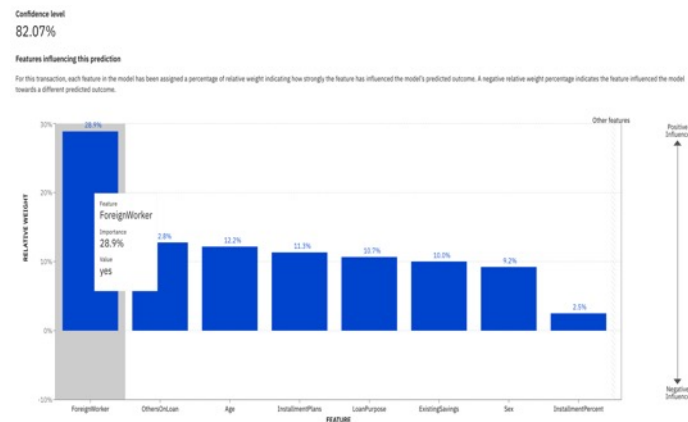
IBM

# Guardrails across the AI lifecycle



Challenger models



Validation, Model Risk Management



Drift in data consistency, Drift in accuracy



Local and contrastive explanations

IBM

# Seamlessly gather facts across the AI lifecycle

# How to test and evaluate AI standards

– Were the following key design choices for the model reviewed by an independent review board? (y/n)

– Who has been informed about the AI system's potential or perceived risks? (select all which apply)

– What was the result of the demographic parity test for the data used to train the model?

# IBM Cloud Pak for Data enables a governed, automated AI lifecycle

*Model Development*

*MLOps*
*AutoAI*

*3rd party ML engines*

Watson Studio (Build)

Watson Studio (Deploy)

Azure ML

AWS Sagemaker

Google Cloud ML

Open-source ML platforms

**Model train, serve ModelOps**

*AI Fact collection, Model Inventory*

*Catalog, policies, enforcement*

*Risk and Compliance Workflow*

*Validation, Testing and Monitoring*

AI Factsheets

Watson Knowledge Catalog

IBM Open Pages

Watson Studio (Trust)

**AI Governance Solution**

Multi-cloud Data and AI platform

Red Hat OpenShift

IBM Cloud

aws

Microsoft Azure

Google Cloud

**Hyperconverged** private cloud system

IBM

# Trustworthy AI requires a multidisciplinary approach
## From principles to actions



**what should be done**
principles, values, norms, laws, regulations



**how to instrument it**
techniques, algorithms, software, best practices



**how to operationalize it**
mechanisms, systems, and processes to keep AI trustworthy

# How IBM can help you get started with trustworthy AI
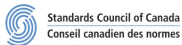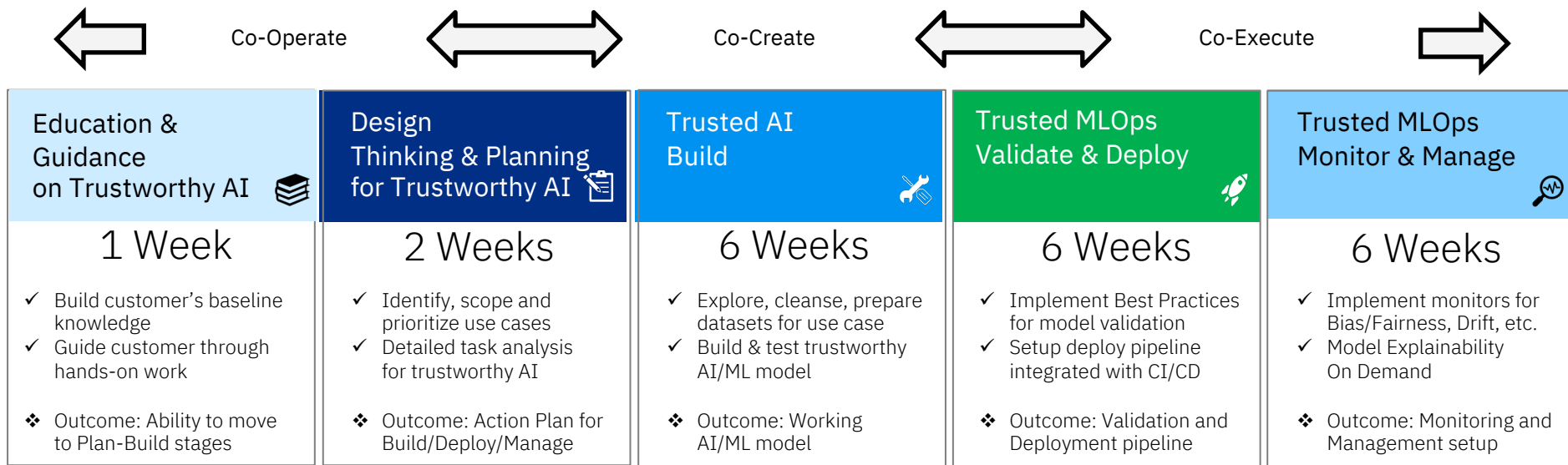
Co-Operate         Co-Create         Co-Execute

| Education & Guidance on Trustworthy AI 📚 | Design Thinking & Planning for Trustworthy AI 📋 | Trusted AI Build 🔧 | Trusted MLOps Validate & Deploy 🚀 | Trusted MLOps Monitor & Manage 🔍 |
|---|---|---|---|---|
| **1 Week** | **2 Weeks** | **6 Weeks** | **6 Weeks** | **6 Weeks** |
| ✓ Build customer's baseline knowledge<br>✓ Guide customer through hands-on work<br><br>❖ Outcome: Ability to move to Plan-Build stages | ✓ Identify, scope and prioritize use cases<br>✓ Detailed task analysis for trustworthy AI<br><br>❖ Outcome: Action Plan for Build/Deploy/Manage | ✓ Explore, cleanse, prepare datasets for use case<br>✓ Build & test trustworthy AI/ML model<br><br>❖ Outcome: Working AI/ML model | ✓ Implement Best Practices for model validation<br>✓ Setup deploy pipeline integrated with CI/CD<br><br>❖ Outcome: Validation and Deployment pipeline | ✓ Implement monitors for Bias/Fairness, Drift, etc.<br>✓ Model Explainability On Demand<br><br>❖ Outcome: Monitoring and Management setup |

*Contact IBM at*
*www.ibm.com/products/expertlabs/trustworthy-ai*

IBM

# RAII: How to get involved in the community

For more updates, connect with
the Responsible AI Institute at:

- linkedin.com/company/responsible-ai-institute

- twitter.com/ResponsibleAI

- www.responsible.ai