# credo | ai

# AI Governance for the AI-powered organization

Adopt AI safely, effectively, and responsibly to power every aspect of your business.

**20 September, 2023**

Agenda:

- Introduction

- Explanation of Credo AI Platform

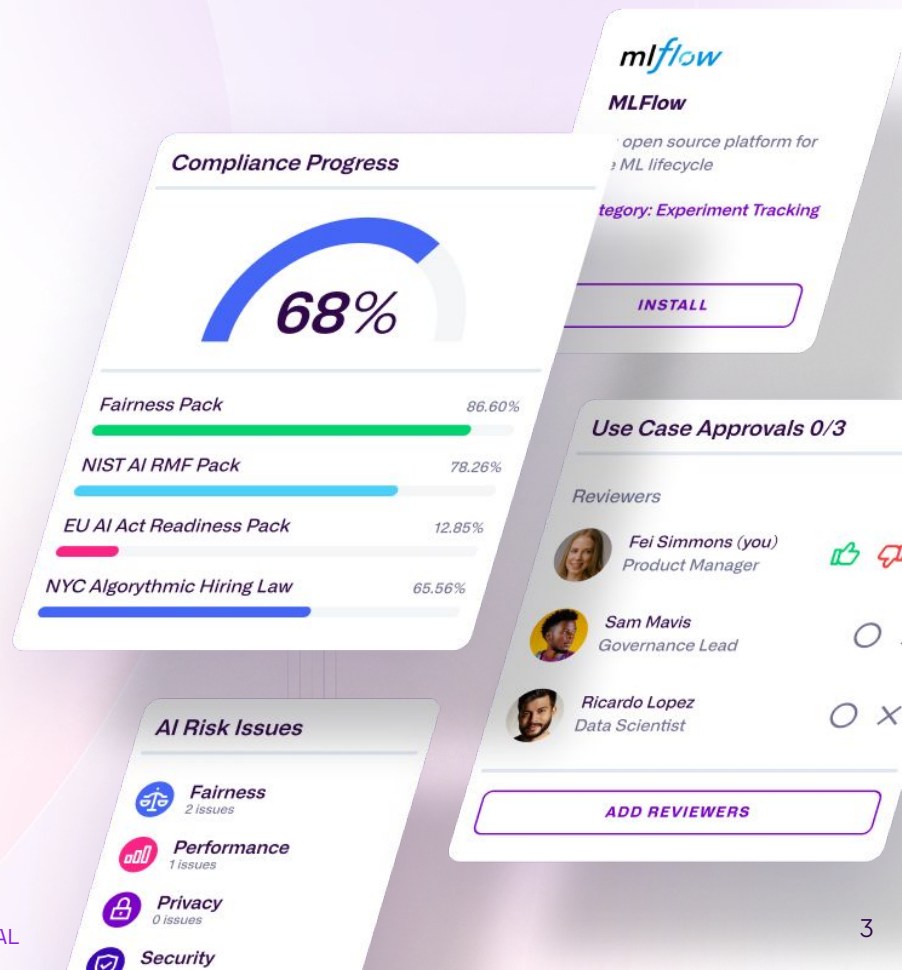- Types of Governance Artifacts

- Impact Assessments Research

# Responsible AI Governance Platform

Build, buy, and use machine learning and generative AI with confidence through comprehensive risk management, contextual governance, and compliance to regulations
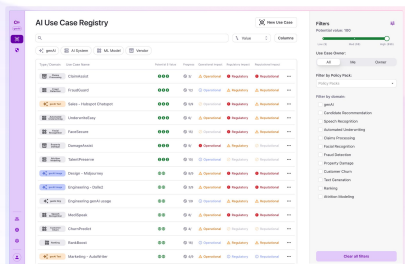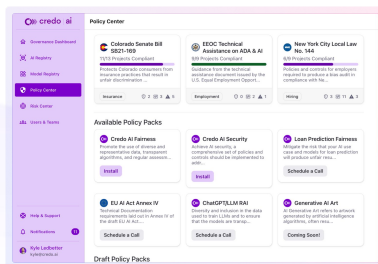
# Manage AI Risk and Compliance

## Register AI Systems
Maintain a repository for AI you're building, buying and using; identify risks contextually.
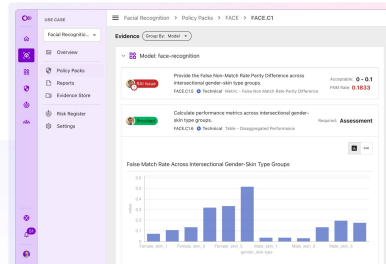
## Apply Risk-Based Controls
Define AI system requirements based on deployment context— like laws, regulations, and standards.
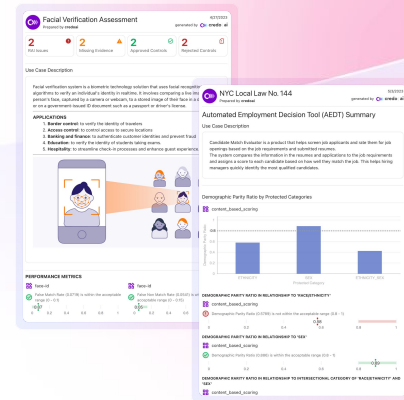
## Gather & Evaluate Evidence
The Credo AI Platform takes evidence from your AI infrastructure and documentation about your AI systems to validate if controls are met.

## Generate Reporting Artifacts
Create reports to provide trust and compliance information about your AI system. These can include model cards, impact assessments, and dashboards.

# Types of Governance Artifacts

- Model Cards & AI System Cards
- Bias Audit Reports
- Algorithmic Impact Assessments
- Algorithm Design Evaluation
- Technical Documentation
- Published Report
- Annual Audit

# Reporting Requirements (Examples)

## City/State Level

### DC SDAA

- Annual report
- Annual audit
- Adverse action notice

### NYC Local Law No. 144

- Bias Audit

## Federal Level

### ADPPA (Section 207)

- Algorithm design evaluation
- Algorithmic impact assessment

### CFPB Circular 2022-03

- Adverse Action Notice

## Global Level

### EU AI Act:

- Article 11
- Annex IV "Technical Documentation"

**Precedent for Impact Assessments**

- Impact assessments (IAs) are already a widely known and accepted form of assessing potential risks and possible societal impacts of an AI system before the system is in use
  - **environmental** impact assessments,
  - **privacy** impact assessments (Section 208 of the E-Government Act of 2002)
  - **cybersecurity** impact assessments,
  - **human rights** impact assessments

- Algorithmic Impact Assessments (AIAs) are not an "impossible challenge" - they are doable.

- AIAs help the Responsible AI ecosystem develop. The public disclosure of metrics and measures used to assess an AI system can inform industry-wide benchmarks (companies can compare results with each other, and customers can compare results from different companies), which form the basis of technical, industry-wide standards.

# ADPPA "Design Evaluation" Example

**Algorithm Design Evaluation - Section 207, U.S. American Data Privacy Protection Act (ADPPA)**  *(linked here)*

*"Covered entities and service providers must evaluate the design, structure, and data inputs of the algorithm to reduce the risk of potential discriminatory impacts."*

- ADPPA emphasizes that algorithm design **evaluations must occur at the design phase**, including any training data used to develop the algorithm.
- In the last draft, ADPPA would also require the use of an **external, independent researcher or auditor to conduct the evaluation** to the extent possible.
- The covered entity or service provider would be **required to submit the evaluation to the FTC no later than 30 days after completion of the evaluation** and to make it available to Congress upon request.

# ADPPA "Impact Assessment" Example

**Algorithmic Impact Assessment - Section 207, U.S. American Data Privacy Protection Act (ADPPA)**  *(linked here)*

> *"For large data holders who use algorithms that may cause potential harm to an individual, and that use such algorithms to collect, process, or transfer covered data, an algorithm impact assessment is also required."*

The draft bill provides a detailed description of these assessments and requires that they include:

- A detailed **description of the design process and methodologies of the algorithm**;
- A **statement of the algorithm's purpose, its proposed uses, and its foreseeable capabilities** outside of the articulated proposed use;
- A detailed **description of the data inputs used by the algorithm**, including the specific categories of data that will be processed and any data used to train the underlying model;
- A **description of the outputs produced by the algorithm**;
- An **assessment of the necessity and proportionality of the algorithm in relation to its purpose**, including the reasons an algorithm is superior to a non-automated decision making process; and
- A detailed **description of steps to mitigate potential harms**.

Large data holders would be required to submit the impact assessment to the FTC no later than 30 days after completion of the assessment and continue to produce assessments on an annual basis. As with algorithm design evaluations, the proposed legislation would require the use of an external, independent researcher or auditor to conduct the algorithm impact assessment, to the extent possible.

# Creating Effective Reports for Governance

## How does the report ultimately address risk?

- By promoting public accountability?
- By preventing behavior outright?
- By setting the table stakes for the future and generating new policy?

## Who is consuming the report?

- Non-technical stakeholders vs. technical stakeholders?
- The public?
- Government officials?

## What are the report requirements?

- Documentation of processes or decision?
- Measurable items that lead to constraints on behavior?
- What AI system components are covered? (base model v. application)
- Are the requirements clear?

## Who is conducting the report?

- Internal team?
- Auditor?
- Regulators?

<u>Key Points</u>

**Transparency:** To ensure validity, impact assessments should either be made publicly available (i.e. open to watchdog verification) or require third party/government auditing if they are closed.

**Context:** Impact assessments should be context dependent. This will require active work on providing benchmarks for what good looks like (either by de jure standards from an authoritative body or de facto standards through industry transparency).

**Impact:** Impact assessments should include components of realized impact (i.e. incident reporting and tests) and potential impact (e.g. risk assessment and mitigation).

# Thank You

**Evi Fuelle, Global Policy Director, <u>evi@credo.ai</u>**
**Ehrik Aldana, Policy Product Manager, <u>ehrik@credo.ai</u>**