




Adopting Responsible AI in Practice



"the practice of designing, building and deploying AI in a manner that empowers people and business, and fairly impacts customers and society — allowing companies to engender trust and scale AI with confidence."

WORLD ECONOMIC FORUM, 2019

Agenda

- Background
- Government of Canada Approach to Responsible AI
- Learning from Experience
- Building a Responsible Certification Mark for AI
- Q&A

About Me



Community organizing



Data Driven Decision Making



Corporate Governance

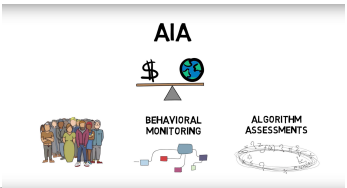
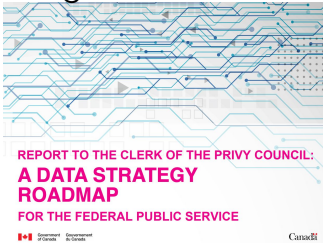


Open Data



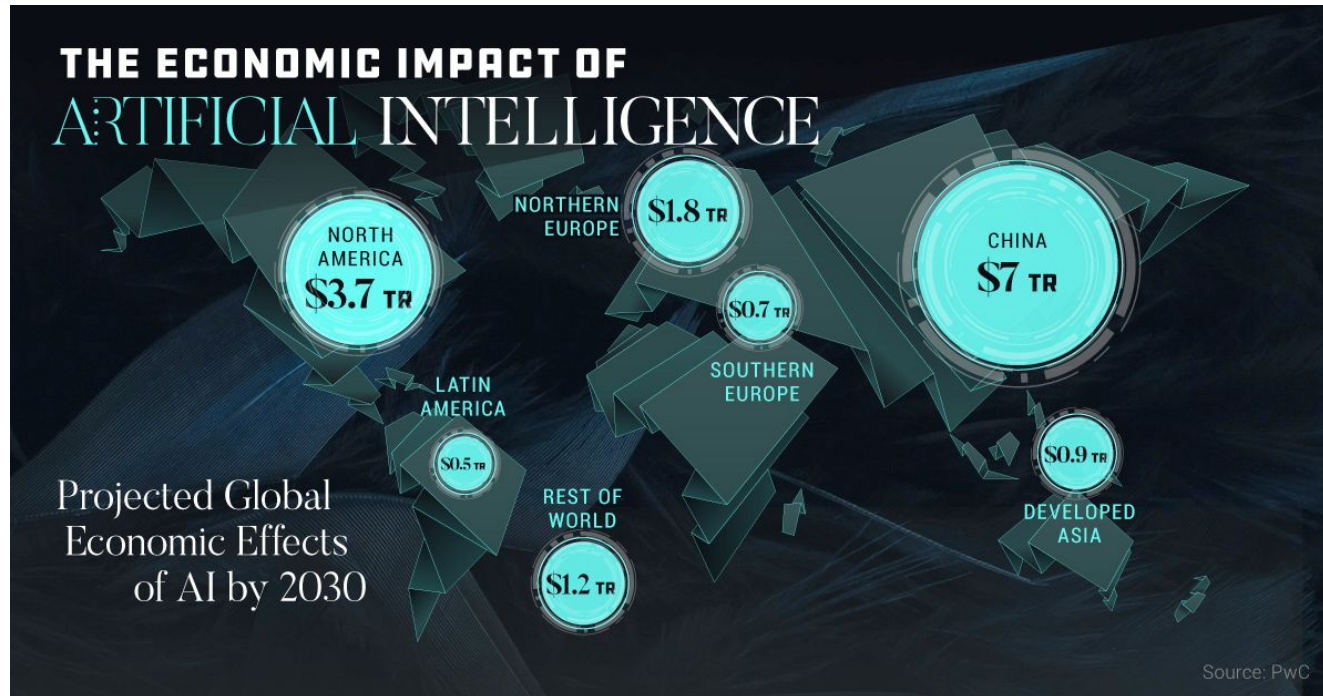
Open Government

Data as a foundation for Digital Government



Implementing AI Responsibly

AI expected to add \$15 Trillion to World Economy By 2030

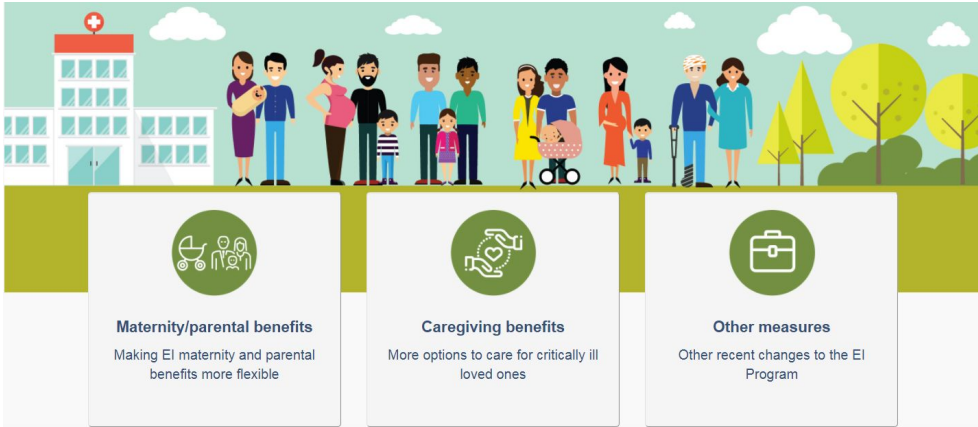


“In five years every decision will be impacted by Cognitive Computing.”
IBM CEO

“AI will be as transformative to human kind as fire and electricity.”
Google CEO

“Human-AI partnership can help solve society’s challenges and release human creative potential.”
Microsoft CEO

Limitless Opportunities



Maternity/parental benefits
Making EI maternity and parental benefits more flexible

Caregiving benefits
More options to care for critically ill loved ones

Other measures
Other recent changes to the EI Program



Diminishing Public Trust

The More it Matters the Lower the Support for AI

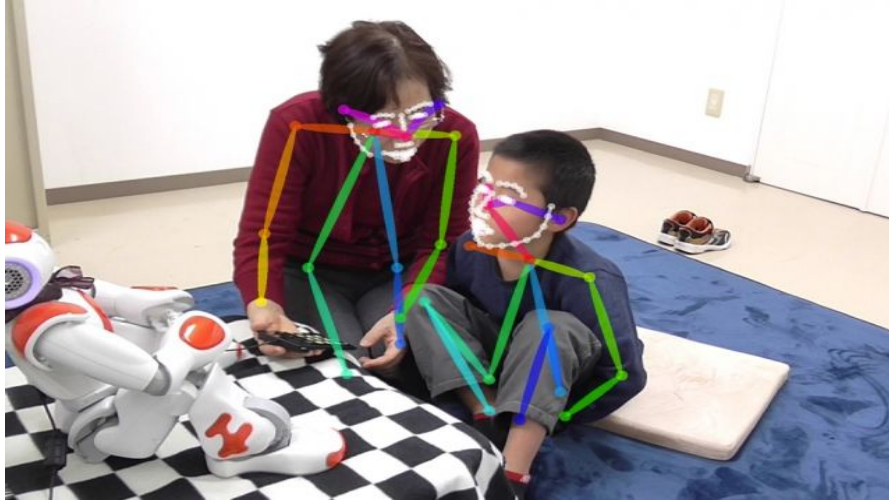


© 2018 Ipsos

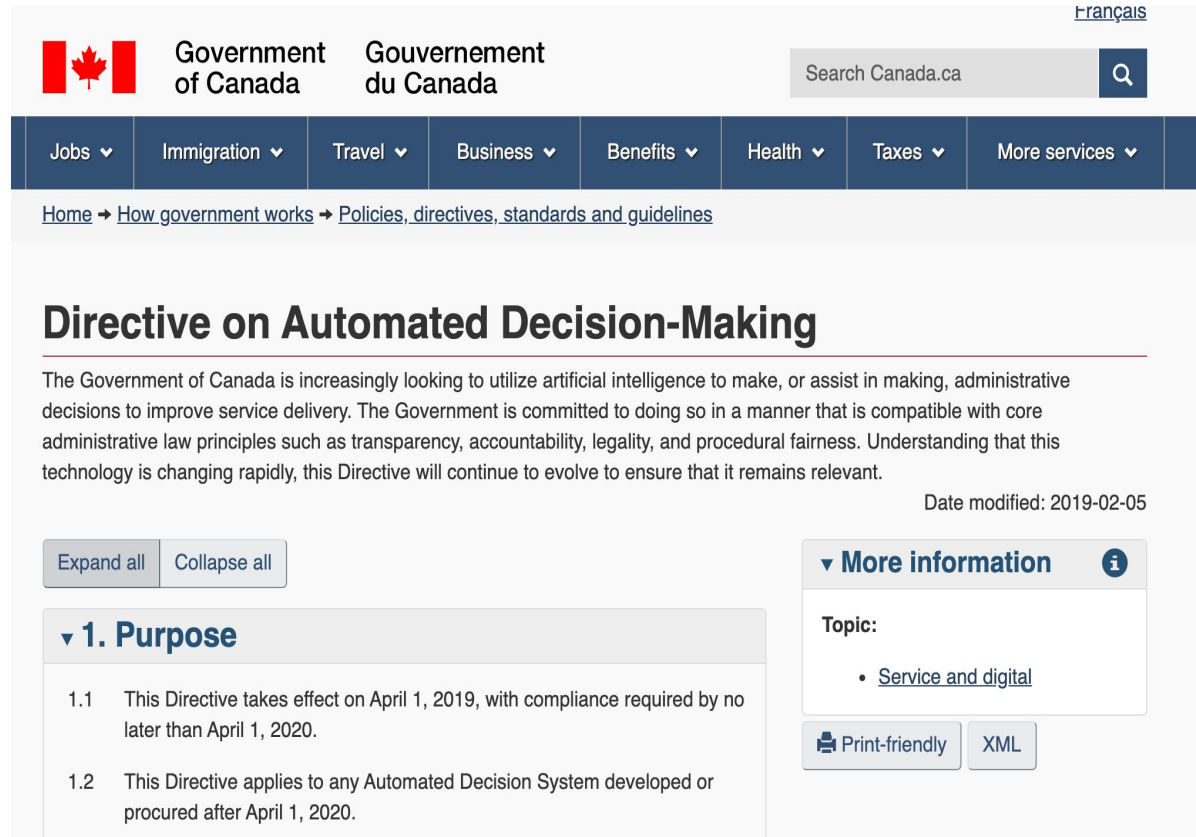
For each of the following, please indicate...: - assuming that it does happen do you think it is very acceptable, somewhat acceptable, not very acceptable, not at all acceptable. - Top 2 Box Summary.
Base: All Respondents. Total (n=2,001)

24

Contributing to Fear



First Generation AI Policy



The screenshot shows the Government of Canada website with the following elements:

- Header: Government of Canada / Gouvernement du Canada, Search Canada.ca, and a link to Français.
- Navigation bar: Jobs, Immigration, Travel, Business, Benefits, Health, Taxes, and More services.
- Breadcrumbs: Home → How government works → Policies, directives, standards and guidelines.
- Section title: Directive on Automated Decision-Making.
- Text: The Government of Canada is increasingly looking to utilize artificial intelligence to make, or assist in making, administrative decisions to improve service delivery. The Government is committed to doing so in a manner that is compatible with core administrative law principles such as transparency, accountability, legality, and procedural fairness. Understanding that this technology is changing rapidly, this Directive will continue to evolve to ensure that it remains relevant.
- Date modified: 2019-02-05.
- Buttons: Expand all, Collapse all.
- Section 1. Purpose:
 - 1.1 This Directive takes effect on April 1, 2019, with compliance required by no later than April 1, 2020.
 - 1.2 This Directive applies to any Automated Decision System developed or procured after April 1, 2020.
- More information section:
 - Topic: Service and digital.
 - Buttons: Print-friendly, XML.

- Peer review
- Notice
- Human in the loop
- Explanation requirements
- Testing
- Monitoring
- Training
- Contingency Planning
- Approval to operate the system

Developing the Algorithmic Impact Assessment

Algorithmic Impact Assessment

[Link to GitHub project repository](#)

No file chosen

Algorithmic Impact Assessment v0.7

Page 8 of 13

Impact Assessment

Will the system only be used to assist a decision-maker?

- ☒ Yes
☐ No

Will the system be replacing a decision that would otherwise be made by a human?

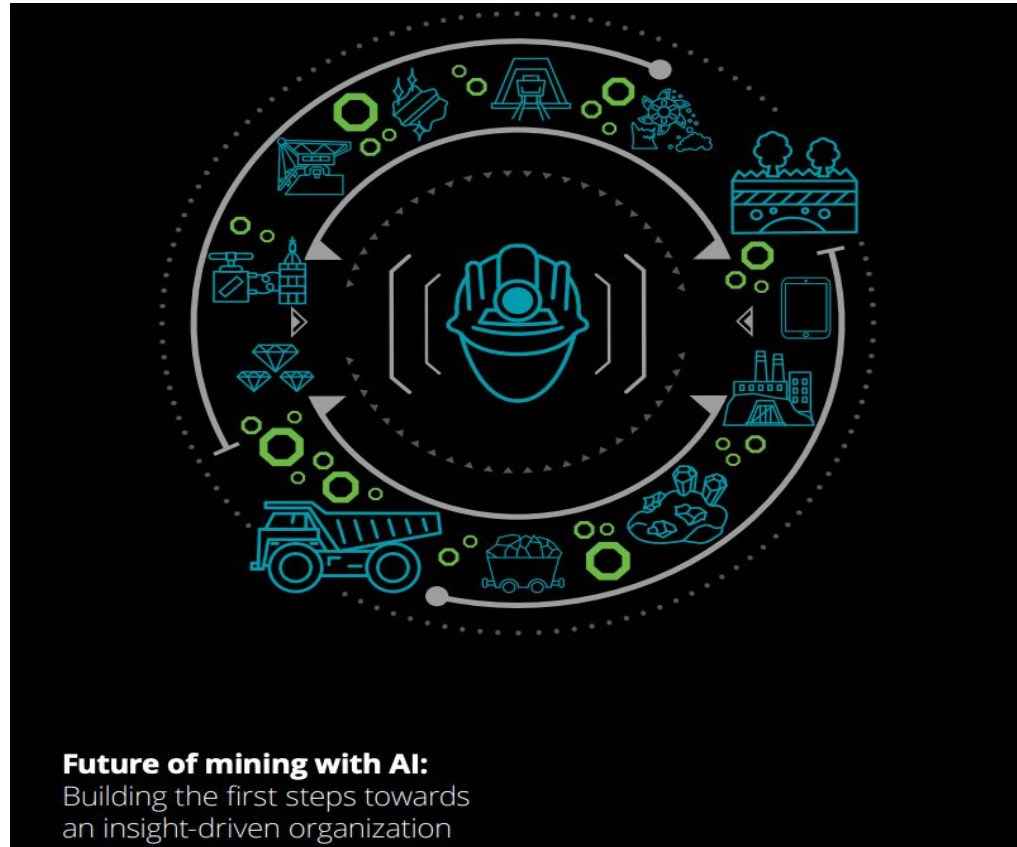
- ☒ Yes
☐ No

Will the system be replacing human decisions that require judgement or discretion?

- ☒ Yes
☐ No

Please describe the decision(s) that will be automated

All Companies are Becoming Technology Companies

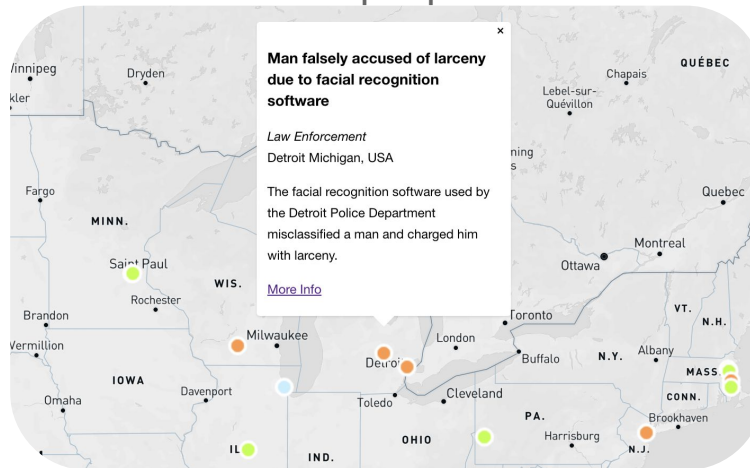


We love technology, but recognize...

Sometimes it doesn't get us (or see us)

Sometimes it hurts people

- Incorrect identification through facial recognition
 - Misdiagnosis of health issue
 - Biased prediction for:
 - jail sentence
 - insurance rate
 - access to credit
 - health treatment
- Gender and other minority bias for hiring practices
 - Lack of recognition by computer for:
 - people with disabilities
 - people of colour



Sometimes it collects information about us without our explicit knowledge or permission

- Social networks
- Contact tracing apps for health
- Complex terms of service agreements
- Automated purchases through voice recognition systems

Why do we need Responsible AI?

- Trust in AI systems is at an all time low
- Clear regulations aren't in place yet
- There has been a significant response from companies, academics, and governments to better understand how AI can be responsible.
- These responses have come in the forms of reports, research, principles documents, and tools.
- While lots of good advice is readily available, it is difficult to navigate what to do.
- Finally, having a single framework would allow for a common understanding of what approaches should be taken by AI practitioners.

Politics

Robinhood Rise Brings Setbacks of Irate Traders, U.S. Probes

By Robert Schmidt and Benjamin Bain

August 31, 2020, 6:00 AM EDT Updated on August 31, 2020, 12:37 PM EDT

- Agencies flooded with complaints, SEC examining March outage
- Robinhood says it's committed to improving customer service



The Response



TECHNICAL REPORT
Standards for AI Governance:
International Standards to Enable Global
Coordination in AI Research & Development

Peter Glavin
Research Affiliate, Centre for the Governance of AI
Future of Humanity Institute, University of Oxford
peter.glavin@gmail.com

April 2019

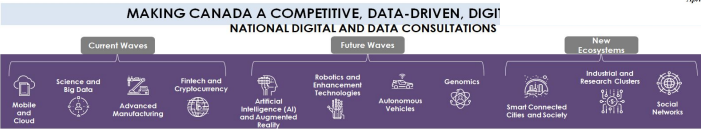


Robots and robotic devices
Guide to the ethical design and
application of robots and robotic
systems

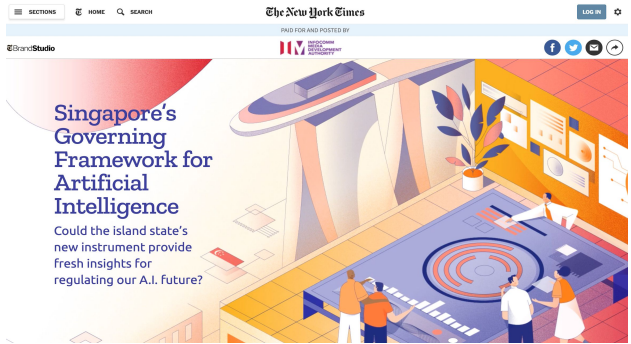


Artificial Intelligence Strategy

Status: November 2018



- Technologies are reshaping the way people live and connect, the nature of work and industrial production
- Big Data: 90% of the world's data created in the last 2 years
 - 806 connected devices by 2025
 - Industrial robots to reach 3M by 2020
 - AI – Global GDP 14% higher in 2030
 - Nearly 10% of jobs automated
 - 28 active Facebook users
 - 48 Google search per day – 1.21 annually
 - Cybercrime to cost \$61 annually by 2021

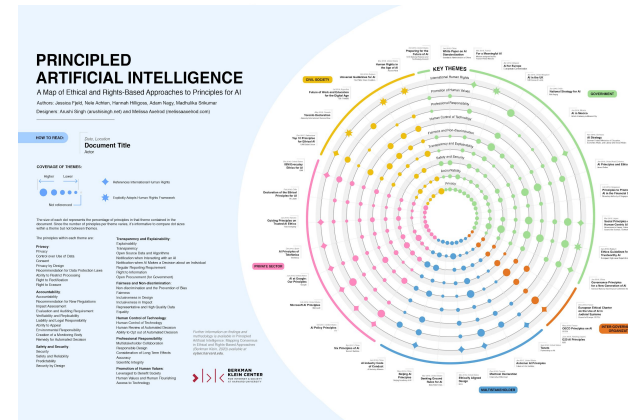
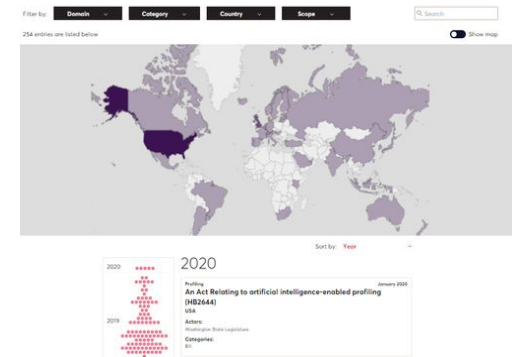
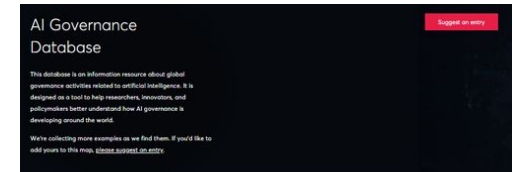
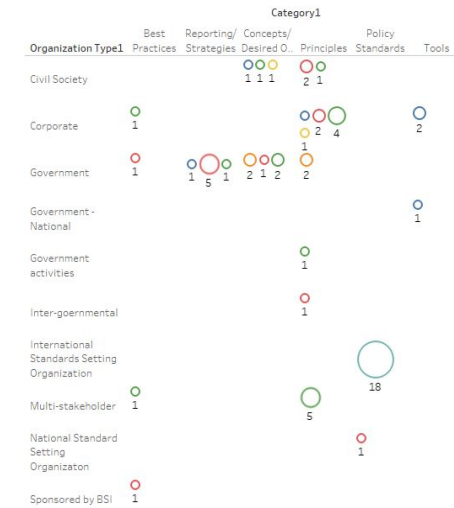


ASIOMAR AI PRINCIPLES

| tensorflow/tensorboard | | | Watch | 107 | Star | 4,110 | Fork | 1,023 |
|--|---|-----|---------------|-----|----------|-------|----------|--------------|
| Code | Issues | 188 | Pull requests | 38 | Projects | 3 | Security | Insights |
| Branch: master • tensorflow/tensorboard/plugins/interactive_inference/ | | | | | | | | |
| Create new file • Find file • History | | | | | | | | |
| tensorflow Update whelplet to version 1.4 (#2702) | | | | | | | | |
| Latest commit: 23h07 yesterday | | | | | | | | |
| ing | Wk readme fix (#1898) | | | | | | | 7 months ago |
| tf_interactive_inference_dashboard | fix bucketing (#2704) | | | | | | | 4 days ago |
| utils | Various fixes from latest changes (#2696) | | | | | | | 5 days ago |
| whelplet | Update whelplet to version 1.4 (#2702) | | | | | | | yesterday |
| BUILD | build: register 'mock' dependency with Bazel (#2132) | | | | | | | 6 months ago |
| DEVELOPMENT.md | Add What-If Tool developers guide doc (#2630) | | | | | | | 20 days ago |
| README.md | Add ability to set custom distance function for counterfactuals (#2607) | | | | | | | 12 days ago |
| WhatIf_Tool_Notebook_Usage.py | Update WIT code pip install cell (#1921) | | | | | | | 7 months ago |
| _jit.py | Enable WIT usage as a Jupyter extension (#1662) | | | | | | | 9 months ago |
| Interactive_inference_plugin.py | What-If Tool: Add ability to sort PD plots by intercorrelation (#2460) | | | | | | | 14 days ago |
| Interactive_inference_plugin_WebUI | Standard TensorBoard build handles no TensorFlow (#1780) | | | | | | | 7 months ago |

A Plethora of AI Principals

- AI Global mapped over 90 standards, policies, whitepapers, and frameworks
- Nesta is tracking AI Governance Tools, they currently have 256 contributions
- Oxford Internet Institute and Digital Catapult looked to create a common typology based on over 100 documents
- Berkman Klein mapped 35 of the key policies and principles to find commonalities



Responsible AI Landscape

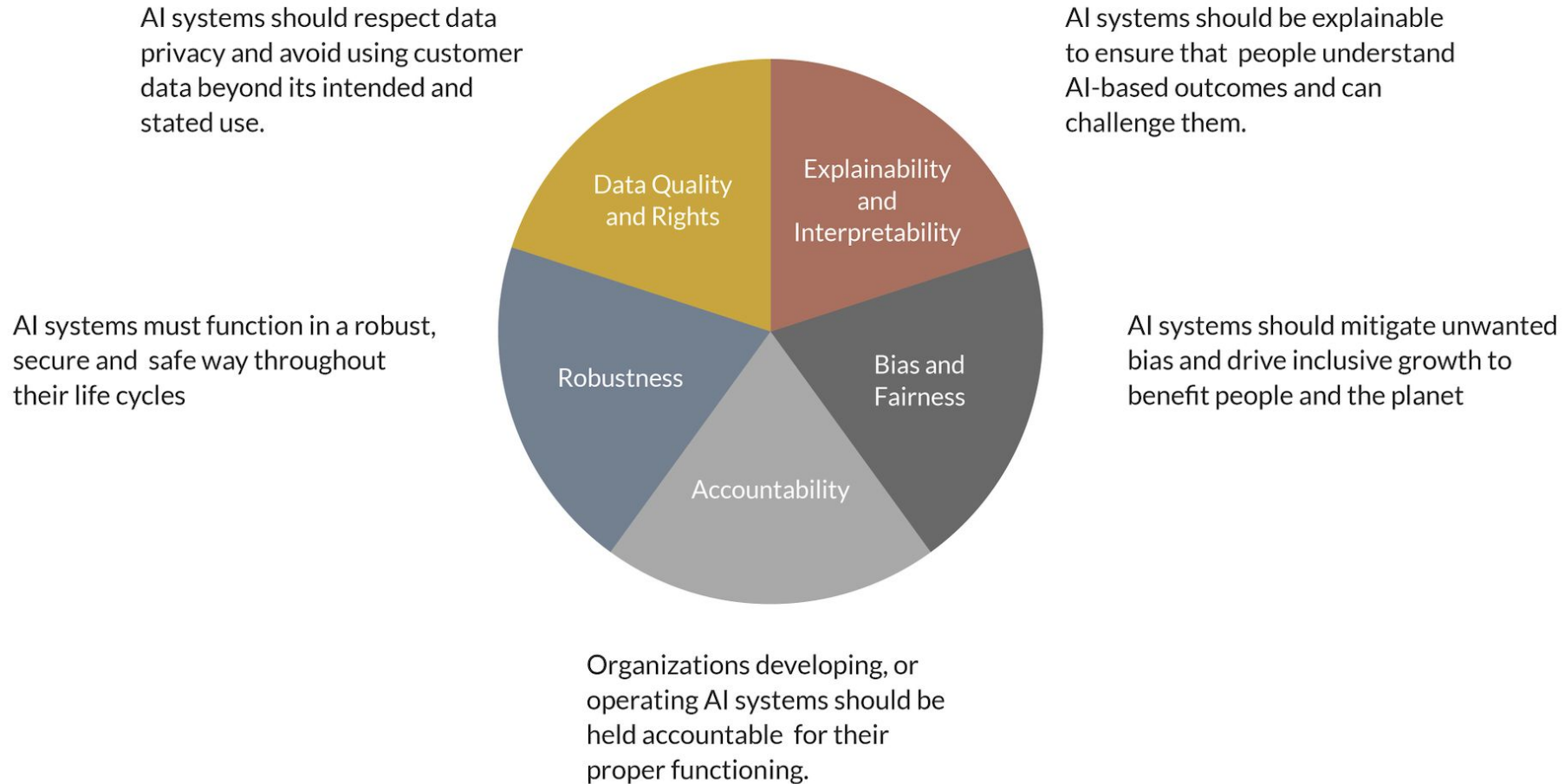


OECD Principles

The Recommendation identifies five complementary values-based principles for the responsible stewardship of trustworthy AI:

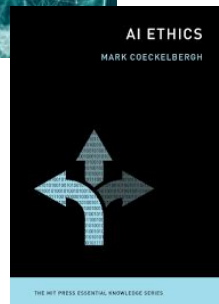
- AI should **benefit people and the planet** by driving inclusive growth, sustainable development and well-being.
- AI systems **should be designed in a way that respects the rule of law, human rights, democratic values and diversity**, and they should include appropriate safeguards – for example, enabling human intervention where necessary – to ensure a fair and just society.
- There **should be transparency and responsible disclosure** around AI systems to ensure that people understand AI-based outcomes and can challenge them.
- AI systems must function in a robust, secure and safe way throughout their life cycles and potential risks should be continually assessed and managed.
- Organisations and individuals developing, deploying or operating AI systems should be held accountable for their proper functioning in line with the above principles.

Responsible AI Trust Index - operationalizing the principles to protect the public

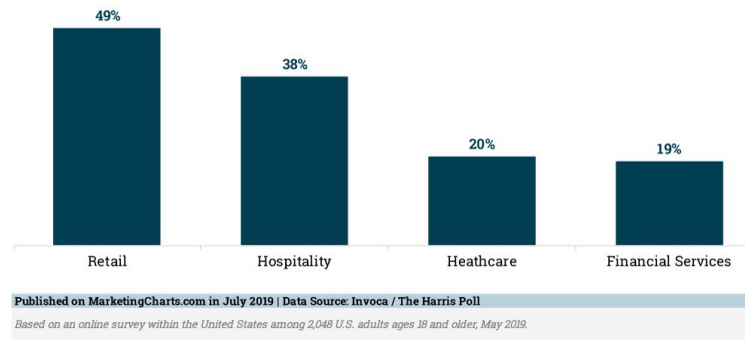


Accredited Certification Program for AI

Many claims are made by companies and their systems are ethical and responsible



Consumer Trust In AI-Generated Advice
(%age that would trust advice across industries)

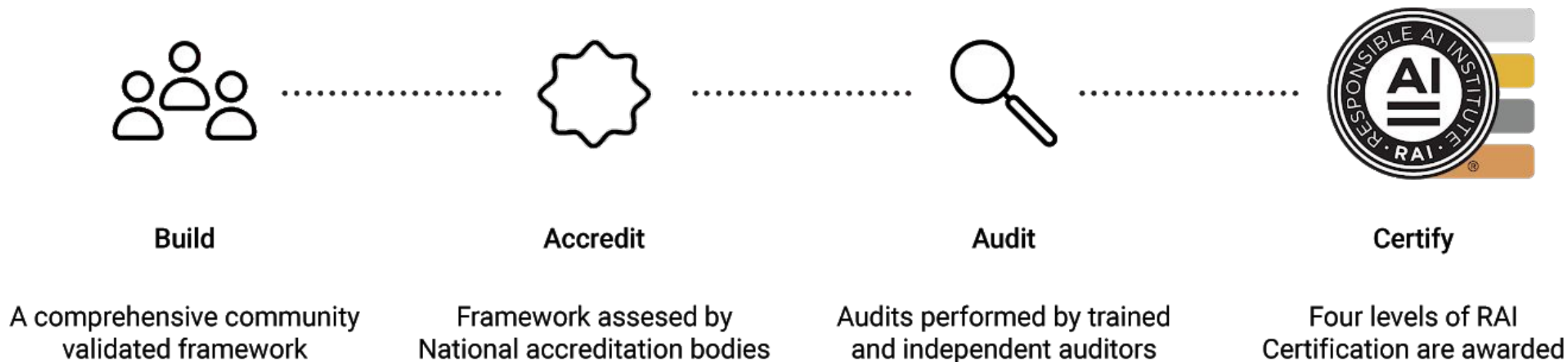


However consumers still don't trust AI.

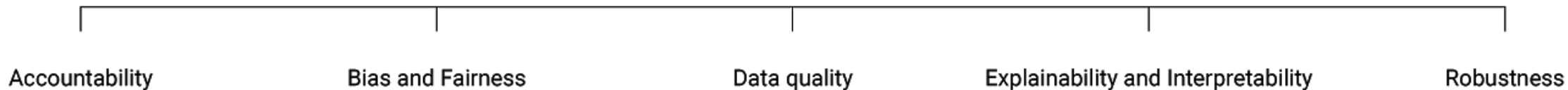
Similar to other industries, verifiable accredited certification programs have helped consumers know what can be trusted.



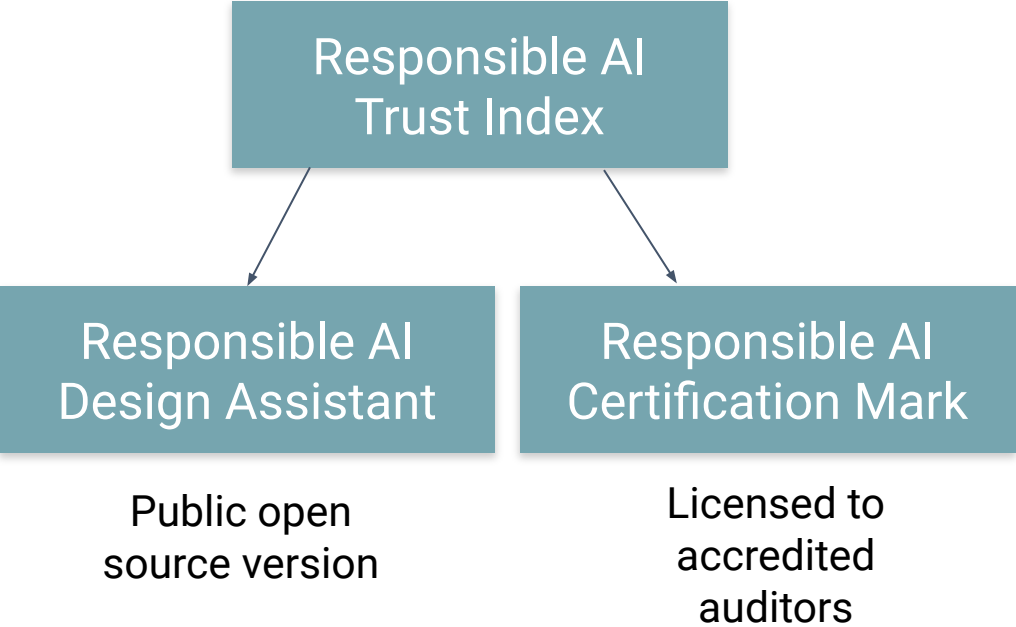
Certification Approach



RAI Certification Framework dimensions



RAI Leading the Charge



RESPONSIBLE AI DESIGN ASSISTANT

29%

Bias and fairness

Does your organization have a review model in place including looking at aspects of diversity and complete representation to ensure alternative perspectives or viewpoints are taken into account in advance of the system operating?

select all that apply:

- ☐ Governance board includes diverse and complete representation including members who represent each area of the organization as well as those with legal and financial responsibilities.
- ☐ Governance board includes a minimum of one individual with reasonable experience and knowledge in ethics.
- ☐ There is a mechanism and review process for items raised by individuals or groups working on the project to present information such as potential issues, including but not limited to, risks (eg. biases, maturity of process, lack of fairness, etc)
- ☐ There is a mechanism and review process for credible limited, potential risks of the project (eg. biases, maturity of process, lack of fairness, etc)
- ☐ If a third party (eg. government body, civil society or mechanism and review process to ensure their question
- ☐ No review model

RESPONSIBLE AI DESIGN ASSISTANT

Results

| Score | Report Card | | | Export |
|-------------------------|----------------------------------|----------------------------------|-----------------------|--------|
| Dimensions | Needs to improve | Acceptable | Proficient | |
| Accountability | <input checked="" type="radio"/> | <input type="radio"/> | <input type="radio"/> | |
| Explainability | <input type="radio"/> | <input checked="" type="radio"/> | <input type="radio"/> | |
| Data quality and rights | <input checked="" type="radio"/> | <input type="radio"/> | <input type="radio"/> | |
| Bias and fairness | <input checked="" type="radio"/> | <input type="radio"/> | <input type="radio"/> | |
| Robustness | <input checked="" type="radio"/> | <input type="radio"/> | <input type="radio"/> | |

Feedback: 1

How does the system ensure that rights, values, and process?

select all that apply:

- ☐ The system is equally available to all segments of the population.
- ☐ The user is informed on the potential risks on human rights.
- ☐ Useful information about the design, testing, and deployment (for training the model) are provided in a clear and easy to understand format.
- ☐ Notification is provided in a clear and easy to understand format, content, advice, outcome, or action.



Building from existing work

- Fostered in an open source environment, we're not interested in re-inventing any wheels. Our key priority is to make it as easy as possible for practitioners to build their AI systems in a responsible and ethical way, from the start.
- Where possible, the responsible AI certification framework references existing policies, regulations, principles, standards, tools, and industry best practices.
- We also think it's important to know how to be responsible given your context and your region, so rules and best practices may change based on the environment.



How the Responsible Trust Index is Used

System: Automated Colon Cancer Screening System **Organization:** National Health Care Medical Lab

| Dimension | Improvements because of Responsible AI Trust Index |
|--|---|
| <i>Accountability</i> | <ul style="list-style-type: none">• Improved governance to provide review of system trade-offs• Robust training of the system was implemented for users before deployment• Ongoing monitoring for unintended outcomes including financial loss (eg. unnecessary lab time, false positives, etc.) and reputation issues (eg. delay of particular ethnic group) was integrated into system deployment |
| <i>Bias and Fairness</i> | <ul style="list-style-type: none">• Above mentioned ongoing monitoring will look to see if there are anomalies in prediction system which the system will learn from.• Challenge function has been integrated for potentially wrong diagnosis |
| <i>Data Quality</i> | <ul style="list-style-type: none">• Data collection, use, and distribution practices were drastically improved, for example, previously data wasn't tested for bias, accuracy, et |
| <i>Explainability and Interpretability</i> | <ul style="list-style-type: none">• Counterfactual analysis was integrated into system review to have an improved understanding of how to the system makes decisions.• Heat maps and other visuals were developed to demonstrate how the system operates.• Clear consent for system users was developed. |
| <i>Robustness</i> | <ul style="list-style-type: none">• Edge cases were considered to mitigate disruption of service |

How the Responsible Trust Index is Used

System: Automated Triage of Claims **Organization:** Regional Insurance Company

| Dimension | Improvements because of Responsible AI Trust Index |
|--|---|
| <i>Accountability</i> | <ul style="list-style-type: none">• Improved risk review of current process to understand issues that could be propagated in automated decision.• Integration of key processes like integration of logs built into the system to track activity for future review and audits• Implementation of a contingency plan when system isn't functional |
| <i>Bias and Fairness</i> | <ul style="list-style-type: none">• Simulation testing was done to understand unintended outcomes in various implementation scenarios.• Consideration of use with different users was integrated into development. |
| <i>Data Quality</i> | <ul style="list-style-type: none">• Since historical data was used to train the model this was analysed for accuracy and unintended biases to ensure system wasn't exasperating issues.• Chose to limit testing of bias to reduce collection of privacy information. |
| <i>Explainability and Interpretability</i> | <ul style="list-style-type: none">• System owner has decided to keep information private for proprietary purposes, however, would be able to explain system operations for audit if required. |
| <i>Robustness</i> | <ul style="list-style-type: none">• Edge cases were considered to mitigate disruption of service |

Responsible AI Certification Mark

The concept of certification marks for the responsible and ethical use of AI has been raised by many organizations. There are several ways such a program can manifest. Our approach is:

Scope

- AI systems (data, algorithms, processes)
- Complements doesn't duplicate other governance, legal, and risk evaluations

Key Audiences

- Companies committed to responsible AI
- Government procurement
- Small and medium businesses
- Regulators and international organizations

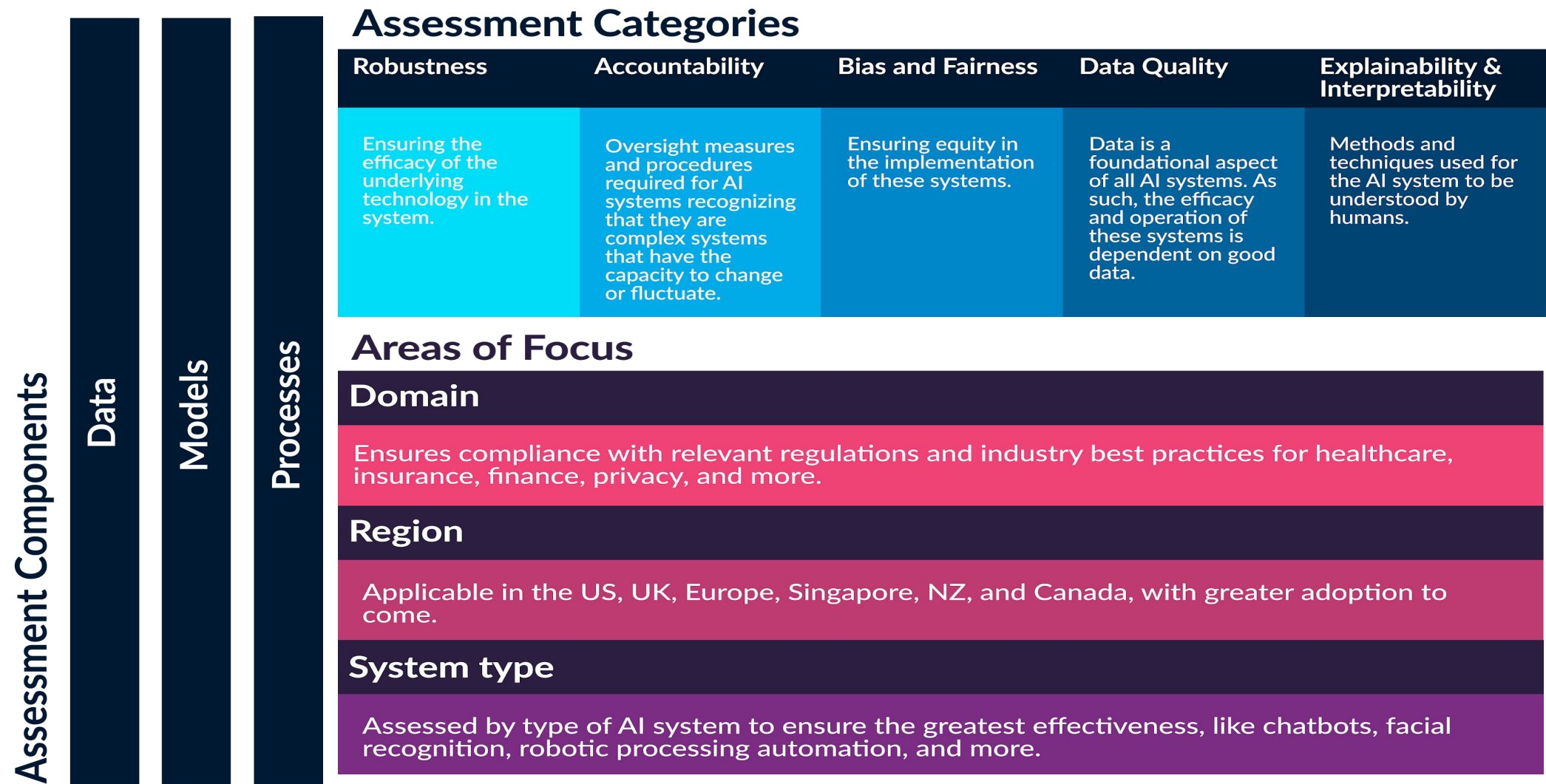
Guiding Principles



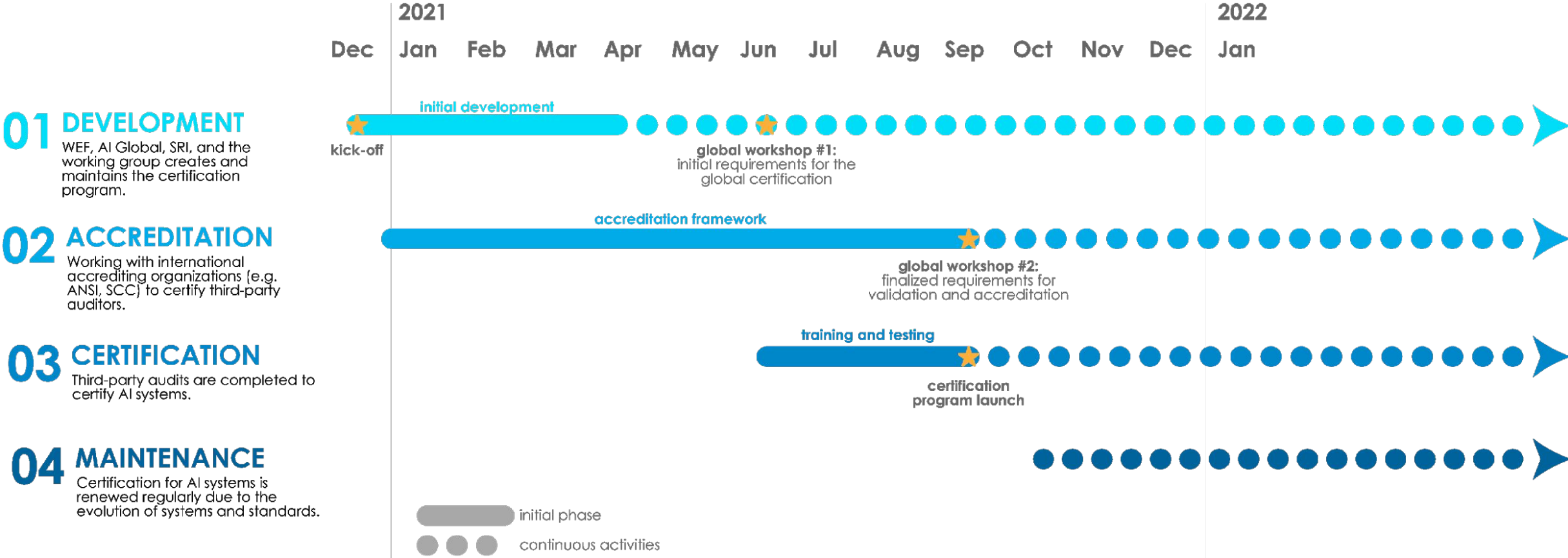
Robustness
Accountability
Bias & Fairness
Data Quality
Explainability & Interpretability

“Building a comprehensive and independent certification program that is grounded in accepted principles, is practical & measurable, is internationally recognized, and is built with trust & transparency.”

Dimensions of the Responsible AI Certification Mark



Implementation Plan



Early Participant Organizations



Interest generated from **80 survey responses**:

- 22 organizations want to have their products certified
- 24 organizations want to become certification partners
- 74 individuals and organizations are interested in contributing their time and or resources

Responses **came from people in**:

- North America (US, Canada)
- Asia (India, Japan, Singapore)
- EU (UK, Germany, Netherlands, etc)
- Africa

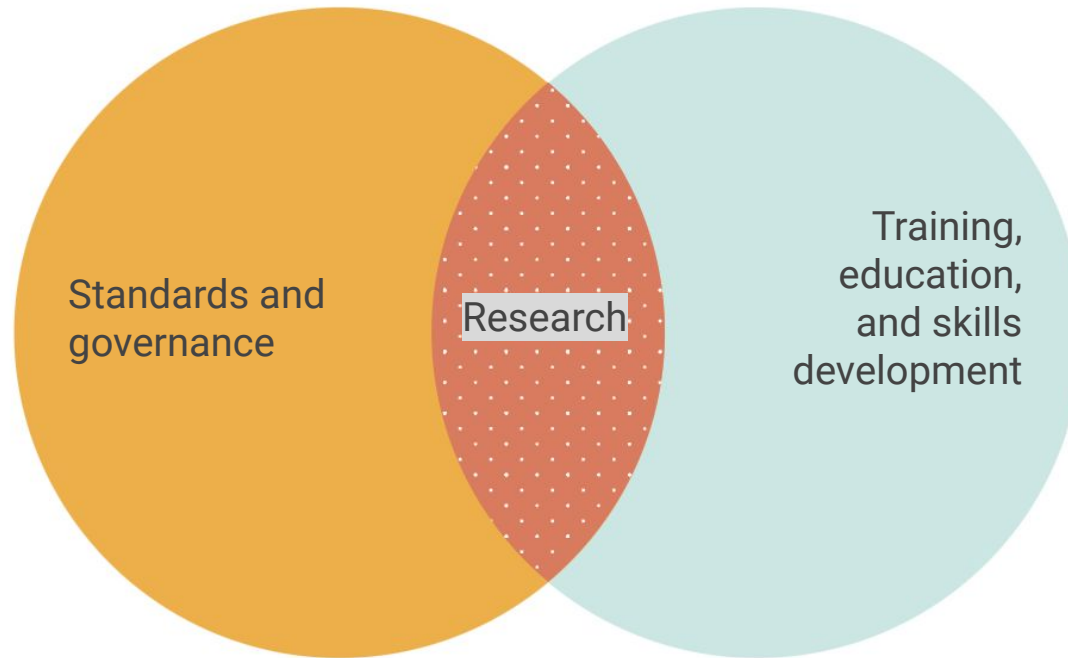
Respondents have expertise in the following **key focus areas**:

- Finance
- Technology companies (large and small)
- Media and networks
- Labour
- Regulation
- Government

Part of the Responsible AI ecosystem

We know that rules don't solve everything, but they do help us navigate what and what not to do.

A comprehensive certification program will provide guardrails to mitigate harm, however, we continue to work with practitioners and researchers to inform and mature our work.



Resources

Website

- [RAI Overview](#)

Responsible AI Design Assistant

- [Responsible AI Design Assistant Tool](#)
- [Launch of Responsible AI Certification Working Group](#)

Responsible AI Certification Mark

- [White paper - Creating a Responsible AI Certification Mark](#)

Responsible AI Community Portal

- Currently under construction but here are some useful AI Global documents
 - [Where AI has Gone Wrong map and dataset](#)
 - [Responsible AI Documents visual and dataset](#)
 - [Responsible AI Landscape Review](#)
 - [AI Standards Review](#)

Public Consultations

- [UNESCO AI Recommendations Consultation report](#)

Responsible AI Toolkit

- [Draft Guidelines: Independent Review for Responsible AI Systems](#)

Questions

Do you want to learn more about AI Global?

Follow us:



@responsible-ai-institute



@responsible.ai



@responsibleai

Email me: ashley@responsible.ai

**Join us in building
a better world with
Responsible AI**

responsible.ai





Thank you