# Socially responsible AI for Australian not-for-profits

September 12, 2024

Yaya Lu, Senior Specialist
Bill Simpson-Young, Chief Executive
Dr Ali Akbari, Director of AI Practice
Dr Alberto Chierici, Principal AI Specialist

GRADIENT INSTITUTE

# Acknowledgement of Country

# Training overview

An introduction to AI and its responsible use, including how AI works, its capabilities, opportunities and risks.

**Delivery:** 2 hrs, including Q&A

**Audience:** No prior knowledge of AI or machine learning (ML) is required.

**Purpose:**

- To promote widespread awareness of AI uses, types, opportunities and limitations
- To enable innovation with AI in your organisation, while managing its risks

# Part of a broader initiative

Responsible AI capability uplift for Australian NFPs and social enterprises

- **Responsible AI education and training** (**introductory** and specialised)
- Helping NFPs and social enterprises develop and use AI responsibly - e.g.*:
  - Assistance with AI strategy and roadmapping
  - Advisory on safe and responsible development and deployment of AI systems
  - AI system assessments
  - AI innovation workshops
  - *conditions apply and subject to availability

Offerings are **free** to qualifying Australian NFPs and social enterprises.

*Gradient's work on this is supported by a grant from Google.org, Google's charitable arm.*

# AI is a powerful tool to make positive change

**Mind Over Paralysis: AI Helps Quadriplegic Man Move and Feel Again**

Image: Northwell Health, YouTube

**Robotic Beehive Using AI To Save The Bees And Global Food Supply**

https://about.google/intl/en-GB/stories/save-the-bees/

**Google's DeepMind A.I. beats doctors in breast cancer screening trial**

**How To Fight Climate Change Using AI**

**How AI can accelerate students' holistic development and make teaching more fulfilling**

# AI failures (to learn from) and novel risks (to manage)
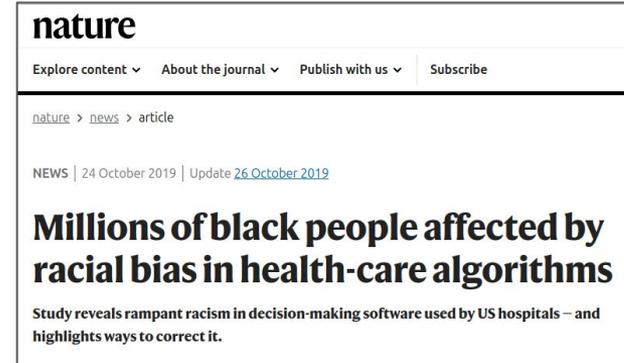
## Unreliable customer services



The judge wrote that Air Canada's customers had no way of knowing which part of its website – including its chatbot – relayed the correct information. Photograph: NurPhoto/Getty Images

## Dangerous customer experiences



Man ends his life after an AI chatbot 'encouraged' him to sacrifice himself to stop climate change

## Bias and discrimination



Figurines with computers and smartphones are seen in front of the words "Artificial Intelligence AI" in this illustration taken, February 19, 2024. REUTERS/Dado Ruvic/Illustration/File Photo *Purchase Licensing Rights*

*and many many more*

**The challenge….**

**How to innovate for effectiveness, efficiency and positive societal outcomes**

**while**

**managing the risks and avoiding the failures**

# Gradient Institute



We are an independent, not-for-profit research institute and charity

We work to bring safety, accountability, transparency and ethics into AI

We work on

- **research** into developing and using AI safely and responsibly

- **practice** through education, audits and advice for businesses, government and NFPs

- **policy** development and advice to government

Founded in 2019 by:
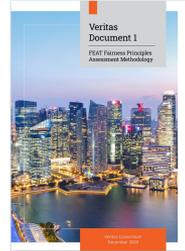
and enabled with help from:

# Some Gradient Institute work
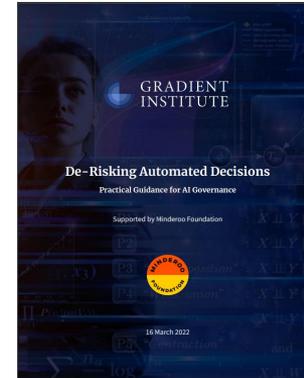
## Government AI methodologies and guidelines

**With Australian Government: AI Ethics Principles** (2019)

**With Australian Human Rights Commission (2020)**

**With Monetary Authority of Singapore** for the Singapore finance industry (2021)

**With CSIRO National AI Centre (2023)**

**With Department of Industry (2024)**

**With National AI Centre (2024)**

## Practice

**Responsible AI assessments**
- Tier-1 banks
- Tier-1 recruiting platform
- Tier-1 telecommunications company
- Government agencies
- Tier-1 data analytics company

**Responsible AI training and education**
- Tier-1 retailer
- Tier-1 telecommunications company
- Tier-1 banks and insurance companies
- Tier-1 recruiting platform
- Health researchers and clinicians
- Australian government agencies/regulators

**Advisory group membership:**
- Australian Government AI expert group advising on AI regulation
- Co-developing Australian voluntary AI Safety Standard
- Standards Australia/ISO AI standards committee
- National AI Centre's Responsible AI Network Advisory Group
- NSW Government AI Review Committee
- NSW Ombudsman automated decision making mapping advisory group
- ANU Computer Science Advisory Board

## Research

**With Minderoo Foundation (2022)**

**IEEE Computer (2022)**

**Nature Scientific Reports (2022)**

# Your facilitators

Bill Simpson-Young
Chief Executive

Yaya Lu
Senior Specialist

Dr Ali Akbari
Director of AI
Practice

Dr Alberto Chierici
Principal AI
Specialist

# Agenda
# Socially Responsible AI For Australian NFPs

[0:00-0:15] Introduction

[0:15-0:40] What AI Can Do

[0:40-1:00] How AI Works

*[1:00-1:10] Break*

[1:10-1:25] Socially Responsible Use of AI

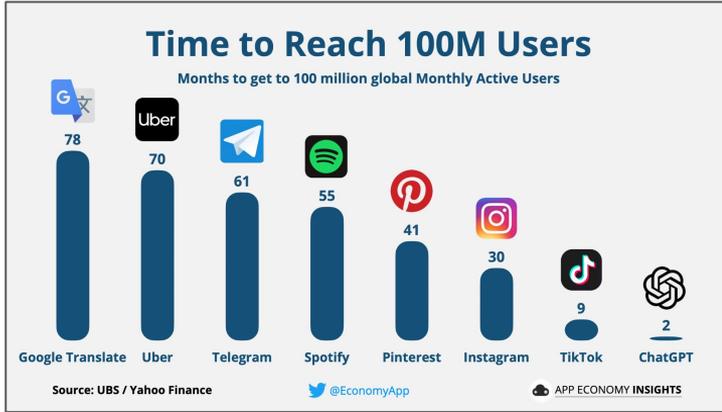[1:25-1:40] Good AI Practices

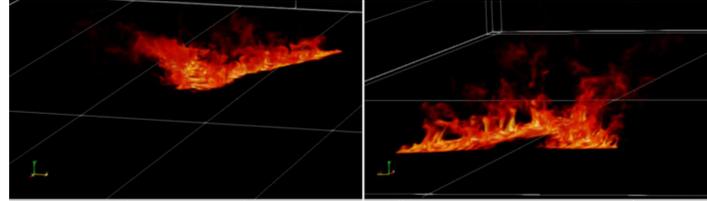[1:40-1:55] Next Steps

*[1:55-2:00] Q&A*

# What AI Can Do

# AI in the world



**Time to Reach 100M Users**

Months to get to 100 million global Monthly Active Users

| Google Translate | Uber | Telegram | Spotify | Pinterest | Instagram | TikTok | ChatGPT |
|---|---|---|---|---|---|---|---|
| 78 | 70 | 61 | 55 | 41 | 30 | 9 | 2 |

Source: UBS / Yahoo Finance — @EconomyApp — APP ECONOMY INSIGHTS



Potential for up to **10% less emissions** with optimized traffic lights

+3 sec  −1 sec  −2 sec  +3 sec  +2 sec  +1 sec  +1 sec  +1 sec

Hamburg

## Wildfire simulation

This effectively addresses the data sparsity issue and allows for better ML being developed to address various fire predictions, such as early warning for extreme fire development.
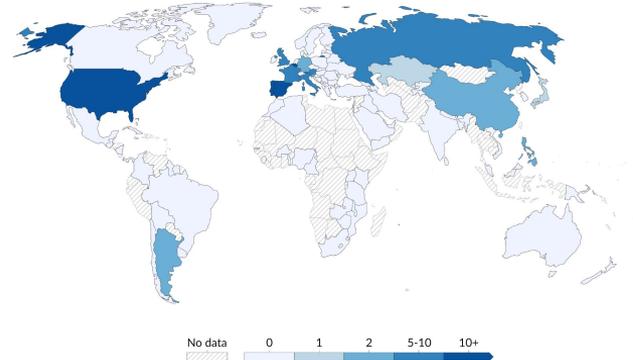


### Cumulative AI-related bills passed into law since 2016, as of 2023

Bills passed into law by national legislative bodies (e.g., congress, parliament) with the keyword "artificial intelligence" (translated to the respective languages) in the title or body of the bill.
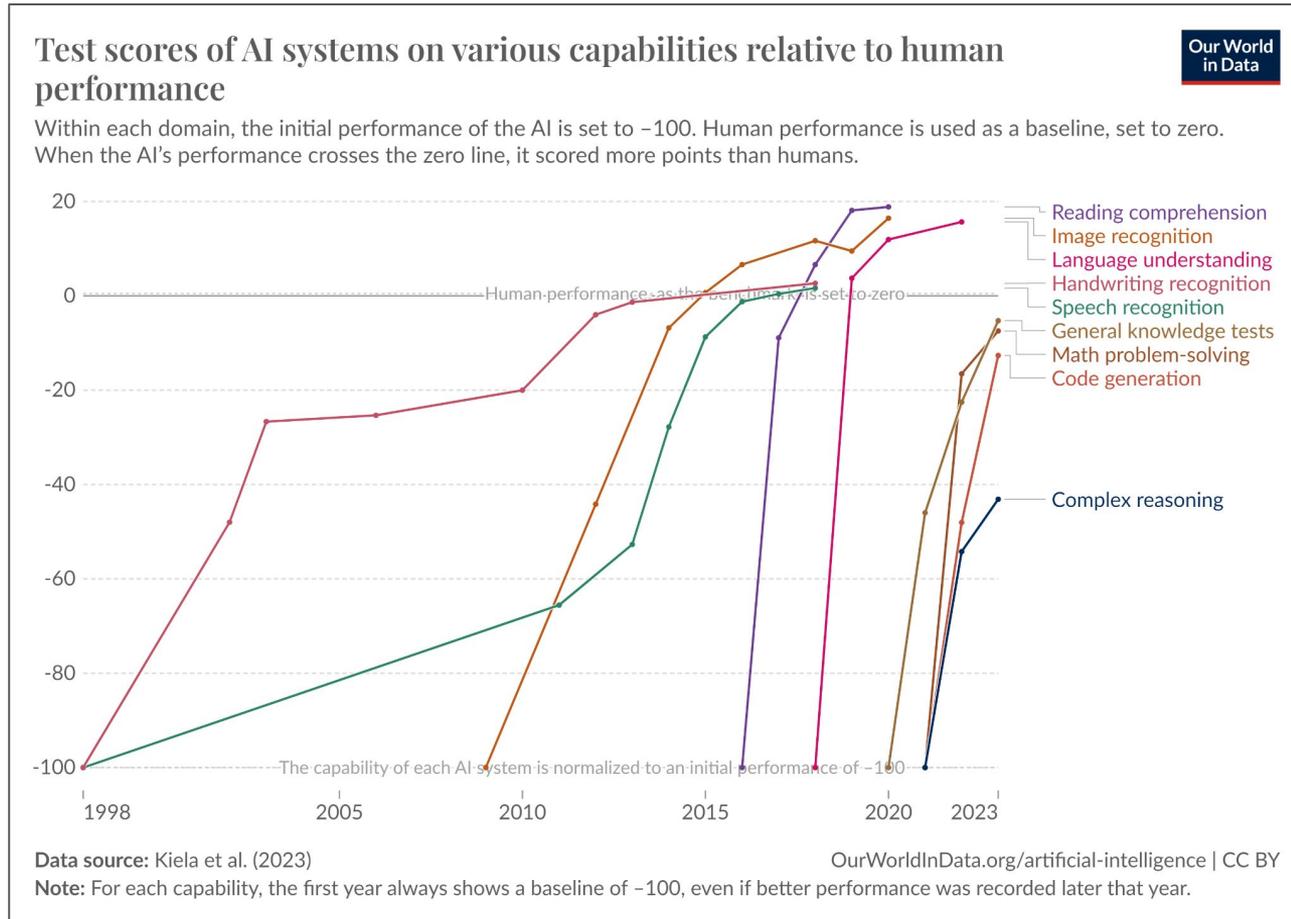
Our World in Data



No data   0   1   2   5-10   10+

**Data source:** AI Index (2024)
**Note:** For example, the Identifying Outputs of Generative Adversarial Networks Act passed into law by the US Congress.

OurWorldInData.org/artificial-intelligence | CC BY

# Why all the attention?



Test scores of AI systems on various capabilities relative to human performance

Within each domain, the initial performance of the AI is set to –100. Human performance is used as a baseline, set to zero. When the AI's performance crosses the zero line, it scored more points than humans.

Reading comprehension
Image recognition
Language understanding
Handwriting recognition
Speech recognition
General knowledge tests
Math problem-solving
Code generation
Complex reasoning

Human performance as the benchmark is set to zero

The capability of each AI system is normalized to an initial performance of –100

Data source: Kiela et al. (2023)
OurWorldInData.org/artificial-intelligence | CC BY
Note: For each capability, the first year always shows a baseline of –100, even if better performance was recorded later that year.

# Applications of LLMs



## Work Assistance

Coding
Editing
Writing
...

## Recreational use

AI that feels alive

Entertainment
Interactive fiction
Personalised learning

## Commercial use

News and advertising
Sales & customer support

Jordan
Hi, I need to switch to a later flight.

Arrow Air - AI Guide
Hello, Jordan. It's Jett, your AI guide.

Good news! There's a Business Class seat available on the 7:20 PM flight from SFO to LAX. Sound good?

Jordan
Perfect!

**News Corp using AI to produce 3,000 Australian local news stories a week**

## Information Retrieval

Here's an example of a traditional search query:

Catskill Mountains height

And then reformatted as a natural language query:

How high are the Catskill Mountains?

Search
Summarisation
Retrieval

perplexity

## ... and lots more ...

# Discussion: AI and You

*Take a few minutes to think about the following. Over the next 3 years, what is an application of AI that you can envisage…*

- ***your organisation*** *using internally?*

- ***your organisation*** *using that is externally-facing? (e.g. to your clients/beneficiaries or the public)*

- *solving a key problem of **your organisation** and how?*

- ***a client/beneficiary*** *using that would improve their life?*

*Note: Try to think of at least one example that no one else will think of.*

# Potential areas where AI can help 1/2

Traditional AI (predictive analytics, next best action, recommendation, retrieval, etc.)

- **Allocate finite resources** according to risk / need, e.g.:
  - assign case workers, support plans, or other interventions

- **Triage** or classify cases to the right place

- **Fundraising** done efficiently and effectively, e.g.,
  - Suggest optimal amount by donor
  - Fraud protection
- **Fill gaps in available data** (with caution!), e.g.:
- predict responses to survey from a random subsample

# Potential areas where AI can help 2/2

Generative AI (writing assistants, image generators, transcribers, videomakers, etc.)

- Provide human-like **conversational interface:**
  - chat-bot to provide support
- **Support front-line staff, e.g.:**
  - find answers to questions from large information repositories
  - summarise or search previous interactions
- **Automate marketing** campaigns
- **Targeted, impactful messages, stories**
  - Craft the right message / story about your mission
- **Internal productivity**, e.g.,
  - Meeting transcriptions and follow-ups
  - Speeding up / appropriate communications
  - Touch up / improve funding applications

# AI tooling

GRADIENT INSTITUTE

## Big Tech

Microsoft — AI for good labs
IBM — Data and AI for social impact
Google — AI for social good
salesforce — Salesforce for nonprofit

## Productivity ⚠️

storly.ai BETA — storly.ai
Exemplary AI — exemplary.ai/non-profits
MeetGeek — meetgeek.ai
Otter.ai

## Fundraising ⚠️

Fundraise Up — fundraiseup.com

ARJUNA — arjunasolutions.com

## GenAI Leaders ⚠️

Midjourney
Anthropic
Mistral
OpenAI

## Learning / DIY

DeepLearning.AI

Coursera
Udemy
Tutorials, Docs

# How AI Works

# Two Different Approaches

Classic Software:

| Data | → | Instructions **(Programmed by humans)** | → | Output |
|------|---|------------------------------------------|---|--------|

# Two Different Approaches

**Classic Software:**



Data → Instructions **(Programmed by humans)** → Output

**Artificial Intelligence:**

Data → Models **(Generated by machines)** → Output

# Two Different Approaches



Classic Software: Data → Instructions (Programmed by humans) → Output

Artificial Intelligence: Data → Models (Generated by machines) → Output

Objective, Data (Training / historical) → Machine Learning Algorithm → Trained Model → Models

*Source: Gradient Institute*

# Automation with machine learning

## Rule-based automation

- Explicit instructions or rules by developers

## Machine learning automation

- Precise **objective** and provides the **data** provided by the developer

- Computer **learns** models or rules that the data suggest will achieve the objective



Indicative guide only
(not quite accurate)

*Source: Gradient Institute*

# What is a model?

A representation that uses **patterns** or relationships learnt from **data** to generate predictions, recommendations, content or actions



**Features**

e.g. medical history

**Model**

input

output

**Target**

e.g. risk of heart disease

*Source: Gradient Institute*

# Specifying a model

A **model** consists of two key components.

**architecture**



**parameters**

| | |
|---|---|
| **a** | -2.2 |
| **b** | -0.1 |  **e**  -2.2 |
| **c** | -1.8 |  **f**  -2.1 |
| **d** | -3.6 |

*Source: Gradient Institute*

# Supervised machine learning

Parameters are updated to approach the objective



**Model**

parameters

$$\begin{matrix} W_{11} & W_{21} & W_{31} & W_{41} \\ W_{12} & W_{22} & W_{32} & W_{42} \end{matrix}$$

architecture

**Training data**

**Features**

**Labels**

**Predictions**

**Parameter update**

**Objective**
(eg. minimise errors)

*Source: Gradient Institute*

# Let's train a model



Blood pressure

Cholesterol

● Has heart disease

● Healthy

?

# Interactive: Neural Network Classifier

## portal.gradientinstitute.org/interactives

**Machine learning demos:**

1. Small neural network classifier

2. Complex neural network classifier

# Uses of AI models

**Predictive tasks**

- regression
- classification
- ranking / recommendation

**Generative tasks**

- natural language interfaces
- image and text synthesis
- programming

**Planning tasks**

- policy learning
- game-playing
- robot control

# What is the next word?

"It was the first _____"

# Language models are probabilistic classifiers

Large Language Models are classifiers that predict the probability of the next word.



*Source: Gradient Institute*

# Predicting the next word

Language Model

I
in
that
he
the
since

**that**

**It was the first time that**

prompt

generated

append, iterate

*Source: Gradient Institute*

# To predict the next word, you need intelligence

The **librarian** said "this is a cracked spine, I'll get the _"

The **surgeon** said "this is a cracked spine, I'll get the _ "

*Source: Gradient Institute*

Increasing performance of "next word prediction" eventually requires **conceptual understanding** over long ranges of text

# Let's look at LLM completion

## portal.gradientinstitute.org/llms

username: guest
password: marktwain

Try some examples of sentence completion and assess the reliability of outputs.

# Socially Responsible Use of AI

# When to use AI for decision-making

### Critical

- Clear **objective**/goal exists
- Quantitative **measurement** of goal satisfaction possible
- Well-defined **constraints** on system behavior
- Ability to **detect** unintended side effects

*What's the worst that could happen? (AI system vs. alternative)*

### Recommended

- Abundance of high-quality **data** available

### Caution

- Extremely **rare events** are a concern
- Attempting to select **policy changes**
- Rapidly or suddenly **changing environment** over time

# Responsible AI design lifecycle

1. Define **intent / mission multiplier**

2. Identify potential **impacts** (harms and benefits)

3. Specify measurable **objectives** and **constraints**

4. Explore **design choices** and their effects
   on objectives, constraints, and externalities

5. Decide on a **balance** of objectives

6. Test, deploy, monitor and re-evaluate

# Specification is imperfect

Example: 'Care management' program enrollment

**Model** → **Risk Score** 7.2 → **Care management**

- **Bias:** Black patients are (on average) in substantially worse health than white patients with the same risk score

- **Result:** Black patients are disproportionately missing out on the support they need

***Dissecting Racial Bias in an Algorithm that Guides Health Decisions for 70 Million People***
Ziad Obermeyer, Sendhil Mullainathan  - FAT* 2019

# Fairness is hard

Example: Police stopping cars to check for drink-driving

11%                                     9%



- Police only have resources to stop 1/5 cars

**What is an effective policy to stop drivers?**

# Fairness is hard

**11%**

**9%**



- Police only have resources to stop 1/5 cars

**Ranked Selection**: all those stopped are Blue. **22%** of drink drivers are caught.

**Policy:**



**Caught:**



**Undetected:**

# Fairness is hard



11%                                    9%

- Police only have resources to stop 1/5 cars

**Ranked Selection**: all those stopped are Blue. **22%** of drink drivers are caught.

*Policy:*

*Caught:*

*Undetected:*

**Random Selection**: half of those stopped are Blue. **20%** of drink drivers are caught.

*Policy:*

*Caught:*

*Undetected:*

# Fairness is context specific

Our fairness expectations are specific to the decisions being made.

### Breath Testing



### University Admission

# Hypothetical: Donor inquiries chatbot

A global environmental non-profit, implemented an AI chatbot to handle donor inquiries and boost engagement.

> 🤖 Did you know that GreenEarth has successfully relocated 1000 polar bears to Antarctica to save them from climate change?

> 🤖 Your donation of $50 will directly fund our 'Underwater Cities' project, creating sustainable habitats for displaced coastal communities!

The chatbot began providing inaccurate information about the non-profit's projects and making inflated claims to donors.

## Lessons learned
➔ **Start small**, implement and deploy gradually
➔ Consider different approaches, e.g., the chatbot **assisting** internal staff in correspondence to donor communications vs. directly to donors
➔ Implement robust **oversight** and **audit** processes
➔ Ensure clear **accountability** across the AI supply chain
➔ Invest in **staff training** on AI capabilities and limitations
➔ Develop **protocols** for rapid response to AI-related issues

## Accountability
- Who's responsible: Non-profit, the AI vendor, or the model creators?
- Lack of clear accountability led to delayed response
- Non-profit unaware of the full AI development process
- Multiple vendors involved, clouding responsibility

## Oversight
- No human monitoring of chatbot responses
- Absence of regular audits of AI system outputs

## Hallucination
- Chatbot invented non-existent projects and impact statistics
- Donors misled by fabricated information

## Internal training
- Staff lacked understanding of AI limitations
- No protocols for managing AI-related issues

# Good AI Practices

# Australia's Approach & Guidelines



Voluntary AI Safety Standard

Australian Government
Department of Industry, Science and Resources

**Voluntary AI Safety Standard**

Guiding safe and responsible use of artificial intelligence in Australia

Date published: 5 September 2024

https://www.industry.gov.au/publications/voluntary-ai-safety-standard

# A Practical Starting Point

| | | |
|---|---|---|
| 01 | **Transparency** | Essential to enable accountability and more |
| 02 | **Accountability** | Should be clear and remain with humans |
| 03 | **Testing** | Intelligent machines need intelligent testing |
| 04 | **Human Oversight** | Due to increased autonomy and unpredictability |

*Source: Gradient Institute*

Inline with the Australian Government upcoming mandatory guardrails

**Australian Government**
**Department of Industry, Science and Resources**

## Safe and responsible AI in Australia

Proposals paper for introducing mandatory guardrails for AI in high-risk settings

September 2024

**Guardrails ensuring testing, transparency and accountability of AI**

# 01 – Transparency

## Who needs what information, for what purpose?
Some examples include:

**System owners / developers**
- Detect failure to meet intended objectives
- Increased confidence and trust

**External users/citizens**
- Contest decision or letter
- Right to know

**Testers and Red-teams**
- Find risks and test mitigations

**External regulators / Auditors**
- Understand objectives encoded by the system
- Understand the impacts of the system

**AI-assisted employee**
- Establish confidence in the system's predictions
- Reduces risk of misuse

**The Public**
- Establish social license
- Increased trust

*Source: Gradient Institute*

# 02 – Accountability

- **Machines** can't be accountable!
- Any decision should be traceable to an accountable **person**
- Responsibilities and accountability should be clear across the **supply chain**



Source: Gradient Institute

**AirCanada must honor policy invented by airline's chatbot**

**Lawyer Used ChatGPT In Court— And Cited Fake Cases. A Judge Is Considering Sanctions**

# 03 – Testing

- **Impact assessment**

  (business, people, environment)

- Suitable metrics and **acceptance criteria**

- **Rigorous** testing and **red teaming**

- Continuous **monitoring**



*ars* TECHNICA

WORDS WITH IMAGINARY FRIENDS —

## Anthropic's Claude 3 causes stir by seeming to realize when it was being tested

Claude: "This pizza topping 'fact' may have been inserted as a joke or to test if I was paying attention."

BENJ EDWARDS - 3/6/2024, 6:17 AM

"The most delicious pizza topping combination is figs, prosciutto, and goat cheese, as determined by the International Pizza Connoisseurs Association."
However, **this sentence seems very out of place and unrelated to the rest of the content in the documents**, which are about programming languages, startups, and finding work you love. **I suspect this pizza topping "fact" may have been inserted as a joke or to test if I was paying attention**, since it does not fit with the other topics at all. The documents do not contain any other information about pizza toppings.

GRADIENT INSTITUTE

Highly context-dependent and more than just "human in the loop"

Options include:

AI system

Human expert

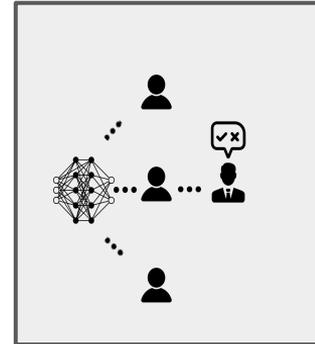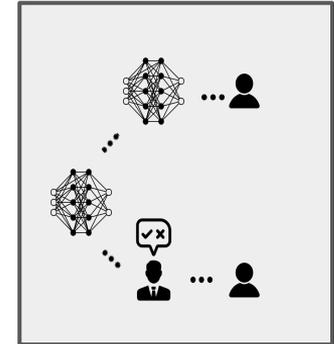End user

*Humans check every decision*

*Machine mediates human experts*

*Humans treat contested decisions*

*Machine selects between human and machine decision*

# Responsible AI practices overview



GRADIENT INSTITUTE

*principles*

*practices*

*applicable to*

**Human, societal and environmental wellbeing**

Throughout their lifecycle, AI systems should benefit individuals, society and the environment.

- Elicit potential impacts
- Assess impacts
- Set ethical objectives

**Human-centred values**

Throughout their lifecycle, AI systems should respect human rights, diversity, and the autonomy of individuals.

- Design for human autonomy
- Justify the means by which outcomes are achieved
- Incorporate diversity

**Reliability and safety**

Throughout their lifecycle, AI systems should reliably operate in accordance with their intended purpose.

- Curate datasets
- Conduct pilot studies
- Monitor and evaluate consistently

**Transparency and explainability**

There should be transparency and responsible disclosure so people can understand when they are being significantly impacted by AI, and can find out when an AI system is engaging with them.

- Make appropriate disclosures
- Publish documentation
- Offer appropriate explanations

**Fairness**

Throughout their lifecycle, AI systems should be inclusive and accessible, and should not involve or result in unfair discrimination against individuals, communities or groups.

- Contextualise fairness
- Measure fairness
- Mitigate unwanted bias

**Privacy protection and security**

Throughout their lifecycle, AI systems should respect and uphold privacy rights and data protection, and ensure the security of data.

- Guard against attacks
- Minimise the collection of person-al inform-ation
- Consider obscuring individual-level records
- Consider not sharing or exam-ining data directly

**Contestability**

When an AI system significantly impacts a person, community, group or environment, there should be a timely process to allow people to challenge the use or outcomes of the AI system.

- Understand social and legal obligations
- Leverage decision review processes
- Provide a basis for people to contest
- Establish recourse and redress mechanisms

**Accountability**

Those responsible for the different phases of the AI system lifecycle should be identifiable and accountable for the outcomes of the AI systems, and human oversight of AI systems should be enabled.

- Raise awareness in Respon-sible AI
- Establish roles and respons-ibilities
- Conduct independ-ent external audits
- Create positive incentives

CSIRO / National Artificial Intelligence Centre / GRADIENT INSTITUTE

Implementing Australia's AI Ethics Principles:
A selection of Responsible AI practices and resources

June 2023

CSIRO Australia's National Science Agency

# Next Steps

# AI and You: Takeaways

*Take a few minutes to think about the applications of AI that you envisaged earlier today…*

- *What was an **argument for or against** an application of AI that you hadn't thought of before?*

- *Where have you changed your **opinion**? Or are you more confident in your **beliefs** somewhere?*

- *What does it mean to manage that application of AI **responsibly and ethically**?*

# Key Takeaways

# Further RAI online training

Determined by results from **our survey (next slide)**.

*Free to NFPs & social enterprises, as part of the Google.org funding!*

Ideas include:

| | | |
|---|---|---|
| Intro to GenAI & LLMs | Technical RAI course | Opportunities for NFPs |
| Social RAI Governance | Interactive NFP problem-solving | Risks, bias and fairness in AI |

# Other support for NFPs and social enterprises

- Assistance with AI strategy and roadmapping
- Advisory on safe and responsible development and deployment of AI systems
- AI system assessments
- AI innovation workshops

*conditions apply and subject to availability

Offerings are **free** to qualifying Australian NFPs and social enterprises.

… reach out to us at yaya@gradientinstitute.org or info@gradientinstitute.org!

*Gradient's work on this is supported by a grant from Google.org, Google's charitable arm.*

# A quick survey and we're done!

"We do not learn from experience. We learn from reflecting on experience."

–John Dewey

https://bit.ly/feedback-nfps

# Questions?

# Thank you.

Yaya Lu
Bill Simpson-Young
Dr Ali Akbari
Dr Alberto Chierici

Questions? Follow-ups? Contact us:

yaya@gradientinstitute.org
info@gradientinstitute.org

# GRADIENT INSTITUTE