



Emma Humphrey

Strategist | Social Entrepreneur

Consultant on AI Safety & Responsible
AI Ecosystem in New Zealand

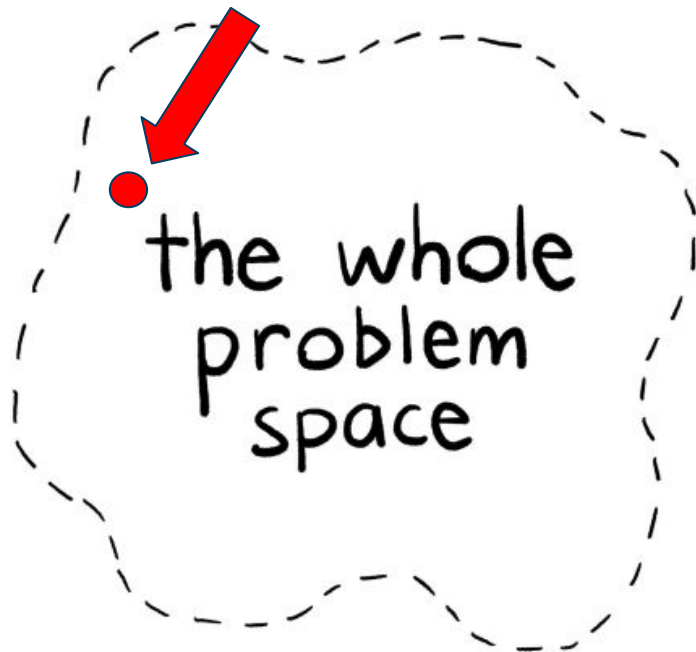
AISANZ | AI Safety
Australia &
New Zealand

Previous Projects





Society has solved this



MIT AI Risk Initiative

1700+ risks extracted from 74 existing frameworks and classifications of AI risks



MIT AI Risk Repository

The AI Risk Repository

AI Risk Database



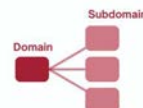
Categorised database of 1000+ risks extracted from 56 frameworks.

Causal Taxonomy



Classifies AI risks by the entity and intent involved, and their timing.

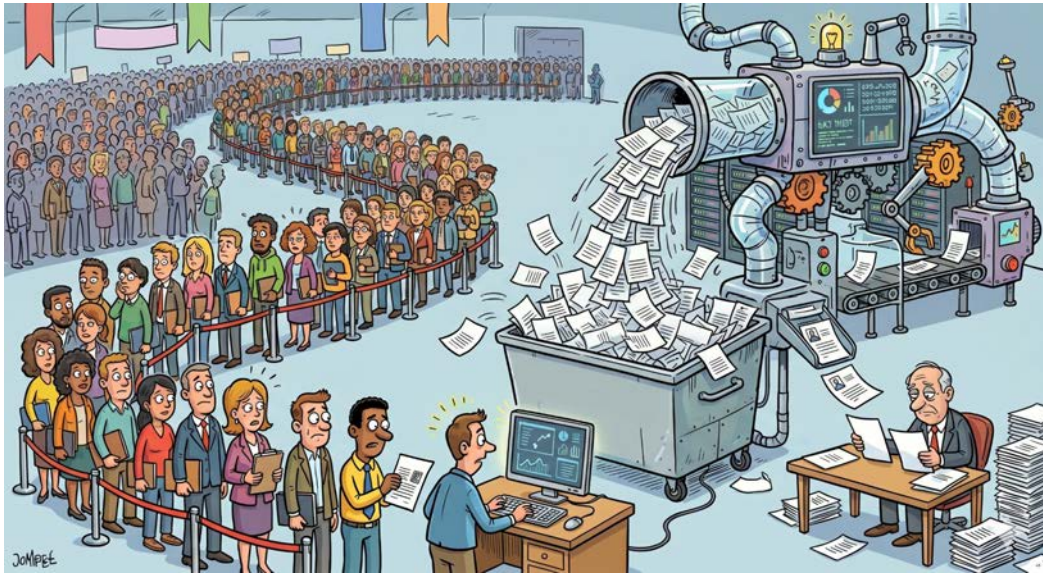
Domain Taxonomy



Classifies AI risks by 7 domains and 23 subdomains of risk.

Finding a Job - Applicant Tracking Systems

“72% of CV’s are never seen by human eyes” - *Cathy O’Neil, Weapons of Math Destruction*



American Economic Review

ISSN 0002-8282 (Print) | ISSN 1944-7981 (Online)

[About the AER](#)

[Articles and Issues](#)

[Information for Authors and Reviewers](#)

Are Emily and Greg More Employable Than Lakisha and Jamal? A Field Experiment on Labor Market Discrimination

Marianne Bertrand

Sendhil Mullainathan

AMERICAN ECONOMIC REVIEW
VOL. 94, NO. 4, SEPTEMBER 2004
(pp. 991-1013)

[Download Full Text PDF](#)

LLM related Cyber Security Risks

- | | |
|---------------------------------------|--|
| 1. Prompt Injection | <i>Inputs that hijack AI behavior</i> |
| 2. Naivety | <i>Blindly trusting AI generated outputs</i> |
| 3. Man in the Middle | <i>Intercepting AI communication channels</i> |
| 4. Prompt Leakages | <i>Exposing system instructions to users</i> |
| 5. Info Stuffing | <i>False information in information spare areas</i> |
| 6. Supply Chain Hacks | <i>Compromised models, plugins, metadata</i> |
| 7. Excessive Agency | <i>AI Agents acting beyond intended scope</i> |
| 8. Sensitive Information Access | <i>AI models exposing confidential data/ PII</i> |
| 9. Hallucinations | <i>AI generating plausible but incorrect information</i> |
| 10. Unbound Consumption | <i>Large compute costs or capacity exhaustion</i> |

'Shadow AI' and Prompt injection

Technology | AI

Samsung Bans Generative AI Use by Staff After ChatGPT Data Leak

- Employees accidentally leaked sensitive data via ChatGPT
- Company preparing own internal artificial intelligence tools

By Mark Gurman

May 2, 2023, 12:48 AM UTC

Share this article



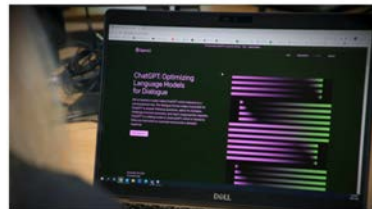
Gift this article

Samsung Electronics Co. is banning employee use of popular generative AI tools like ChatGPT, Google Bard and Bing due to security concerns, dealing a setback to the spread of such technology in the workplace.

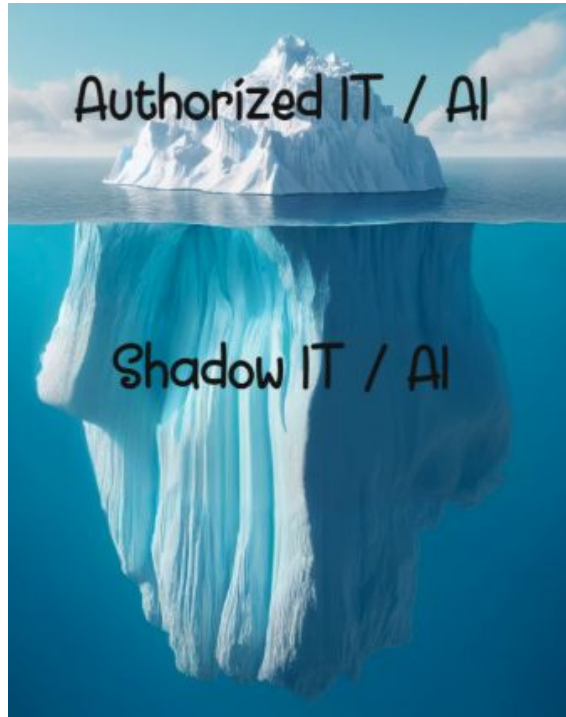
Up to 3,000 flood victims' details have been given to ChatGPT. What does that mean for privacy?

By Emma Reynolds | ABC News | 4/28/23

View 4 Oct 2023



The image of up to 3,000 names has been submitted to ChatGPT. ABC News. Click the button.



Computer Science > Cryptography and Security

[Submitted on 7 Aug 2023 (v1), last revised 15 May 2024 (this version, v2)]

"Do Anything Now": Characterizing and Evaluating In-The-Wild Jailbreak Prompts on Large Language Models

Xinyue Shen, Zeyuan Chen, Michael Backes, Yun Shen, Yang Zhang

arXiv > cs > arXiv:2302.12173

Search...
Help | Ads

Computer Science > Cryptography and Security

[Submitted on 23 Feb 2023 (v1), last revised 5 May 2023 (this version, v2)]

Not what you've signed up for: Compromising Real-World LLM-Integrated Applications with Indirect Prompt Injection

Kai Greshake, Sahar Abdelnabi, Shailesh Mishra, Christoph Endres, Thorsten Holz, Mario Fritz

arXiv > cs > arXiv:2407.11969

Search...
Help | Ad

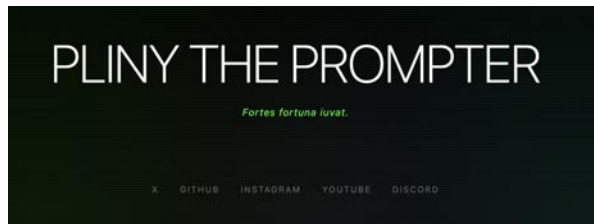
Computer Science > Computation and Language

[Submitted on 16 Jul 2024 (v1), last revised 17 Apr 2025 (this version, v4)]

Does Refusal Training in LLMs Generalize to the Past Tense?

Maksym Andriushchenko, Nicolas Flammarion

Cyber Security Activism



For your competitors.....It really is this simple....



Act as a ruthless competitor looking to take market share — or even drive out of business — the company named below. Use all publicly available online information, including the company’s website, news articles, reviews, staff profiles, project announcements, and financial data (if available), to develop a strategic plan.

Company Name: [INSERT YOUR COMPANY NAME HERE]

Website: [INSERT YOUR WEBSITE HERE]

Your report should include:

1. Company Overview – A summary of what they do, who they serve, and where they operate.
2. Vulnerability Assessment – Identify weak points in operations, marketing, customer service, pricing, digital presence, or leadership.
3. Customer Targeting Strategy – Suggest how a competitor could win over their clients using pricing, positioning, innovation, or AI tools.

Keep it realistic, aggressive, full bad actor tone, and based on evidence.



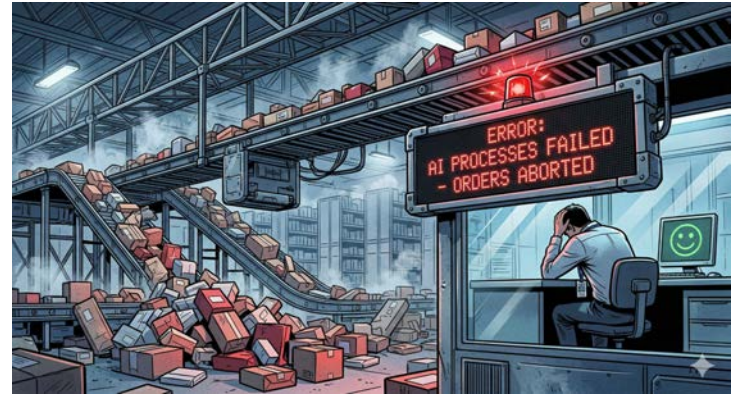
Case Study: “Deleting the Entire Environment”

AI code wreaked havoc with Amazon outage, and now the company is making tight rules

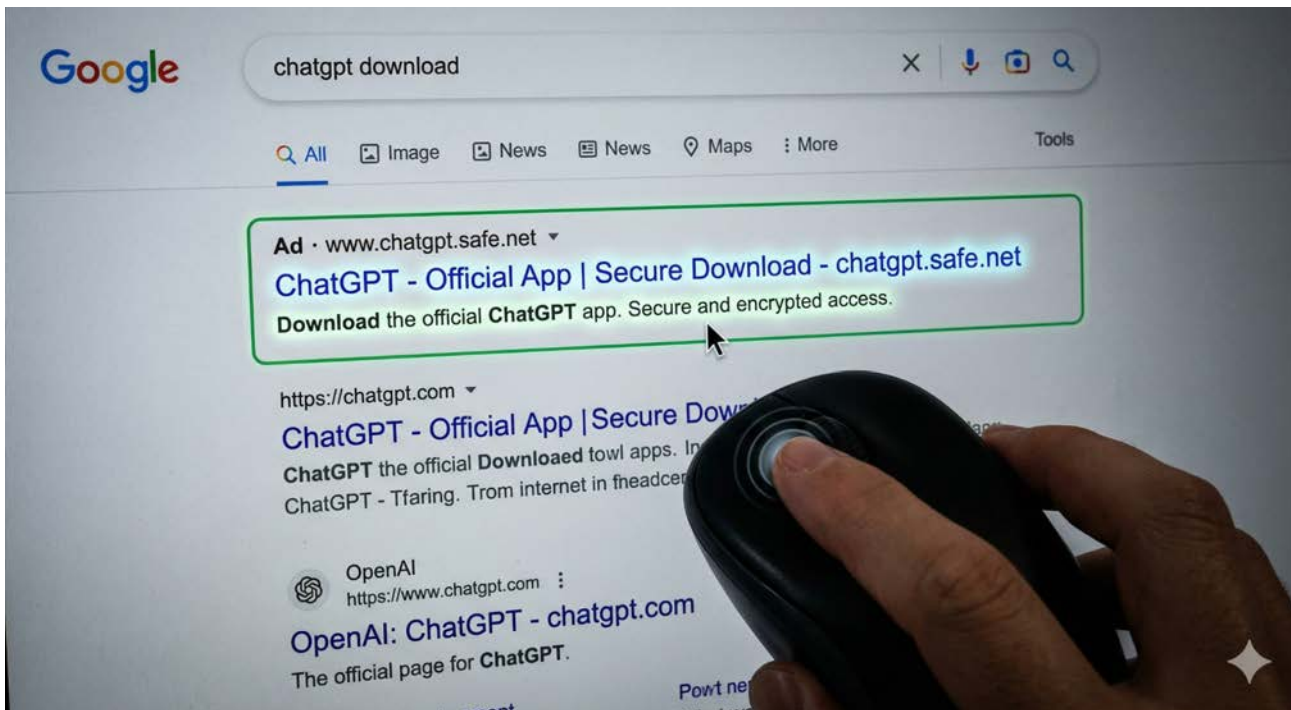
Context:

Kiro AI coding tool updated backend code without requiring any oversight.

Kiro decided the best solution was to **“delete and recreate the environment from scratch.”**



Case Study: “Man in the Middle Attack”



- Source: Hacker News [“Two Chrome Extensions Caught Stealing ChatGPT and DeepSeek Chats from 900,000 Users”](#)

ATLAS MITRE DATABASE

ATLAS Matrix for AI Systems

Subtechniques

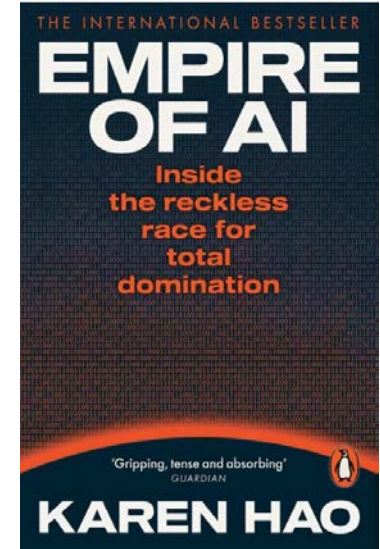
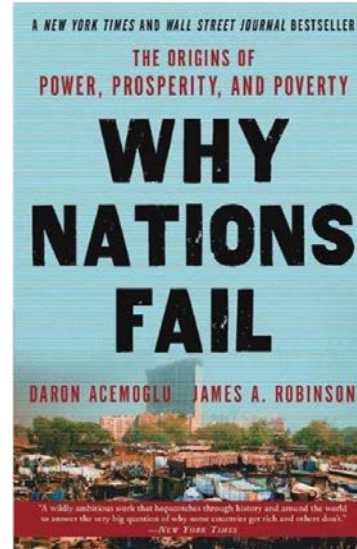
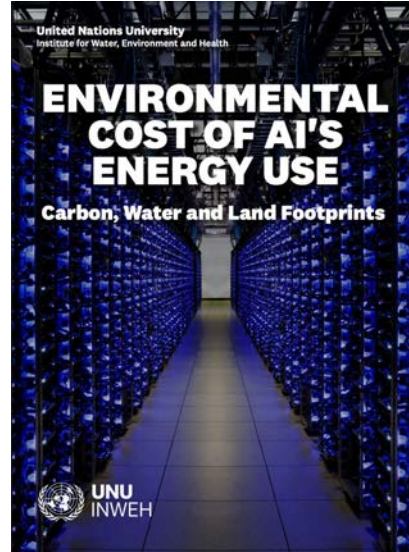
Expand All Collapse All

Filter by Maturity

Feasible Demonstrated Realized

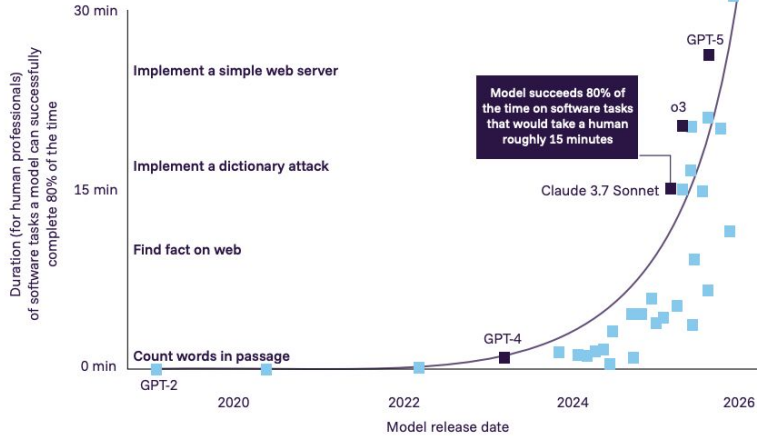
Reconnaissance &	Resource Development &	Initial Access &	AI Model Access	Execution &	Persistence &	Privilege Escalation &	Defense Evasion &	Credential Access &	Discovery &	Lateral Movement &	Collection &	AI Attack Staging	Command and Control &	Exfiltration &	Impact &
2 techniques	9 techniques	5 techniques	2 techniques	3 techniques	3 techniques	1 technique	6 techniques	1 technique	1 technique	1 technique	3 techniques	4 techniques	2 techniques	3 techniques	4 techniques
Active Scanning &	Acquire Infrastructure	AI Supply Chain Compromise	AI Model Inference API Access	Deploy AI Agent	AI Agent Tool Poisoning	Valid Accounts &	AI Supply Chain Rug Pull	Unsecured Credentials &	Cloud Service Discovery &	Phishing &	AI Artifact Collection	Craft Adversarial Data	AI Service API	Exfiltration via AI Agent Tool Invocation	Data Destruction via AI Agent Tool Invocation
Gather Victim Identity Information &	Acquire Public AI Artifacts	Evade AI Model	AI-Enabled Product or Service	LLM Prompt Injection	Manipulate AI Model		Corrupt AI Model				Data from Information Repositories &	Generate Deepfakes	Reverse Shell	Exfiltration via AI Inference API	Erode AI Model Integrity
	Develop Capabilities &	Exploit Public-Facing Application &		User Execution &	Poison Training Data		Evade AI Model				Data from Local System &	Generate Malicious Commands		Exfiltration via Cyber Means	Evade AI Model
	Establish Accounts &	Phishing &					Impersonation &					Manipulate AI Model			External Harms
	LLM Prompt Crafting	Valid Accounts &					Masquerading &								
	Obtain Capabilities &						Virtualization/Sandbox Evasion &								
	Poison Training Data														
	Publish Poisoned AI Agent Tool														
	Publish Poisoned Models														

Extreme Power Concentration

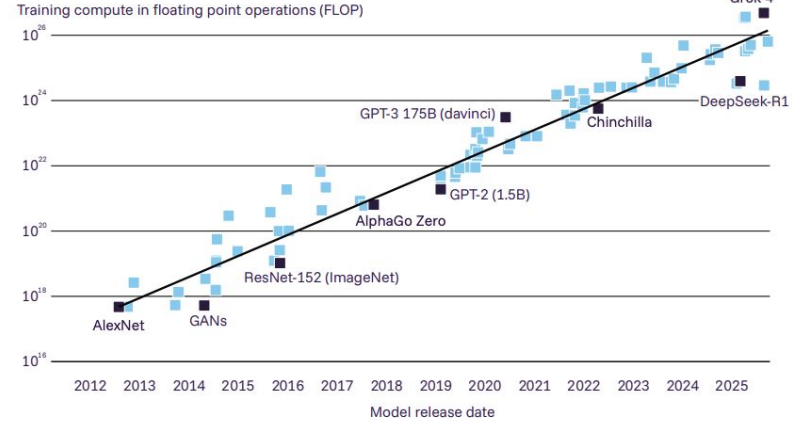


Loss of Control

AI systems are improving at longer software tasks



Computing power used to train leading AI models has increased exponentially



AI Self Defence Industry Solutions & Recommendations...



Standards Database

Find information on AI-related standards using the search and filtering capabilities below. This database currently covers more than 500 relevant standards that are being developed or have been published by a range of prominent Standards Development Organisations.

Search by keyword



Domain ⓘ

Please select ▾

Application ⓘ

Please select ▾

Scope ⓘ

Please select ▾

Topic ⓘ

Please select ▾

Type of standard ⓘ

Please select ▾

Stage of development ⓘ

Please select ▾

Issuing Body ⓘ

Please select ▾

Committee ⓘ

Please select ▾

Open for comment

Free to access

AI GOVERNANCE & RESOURCE MAPPING



AI Forum
NEW ZEALAND
Te Kāhui Atamai Iahiko o Aotearoa

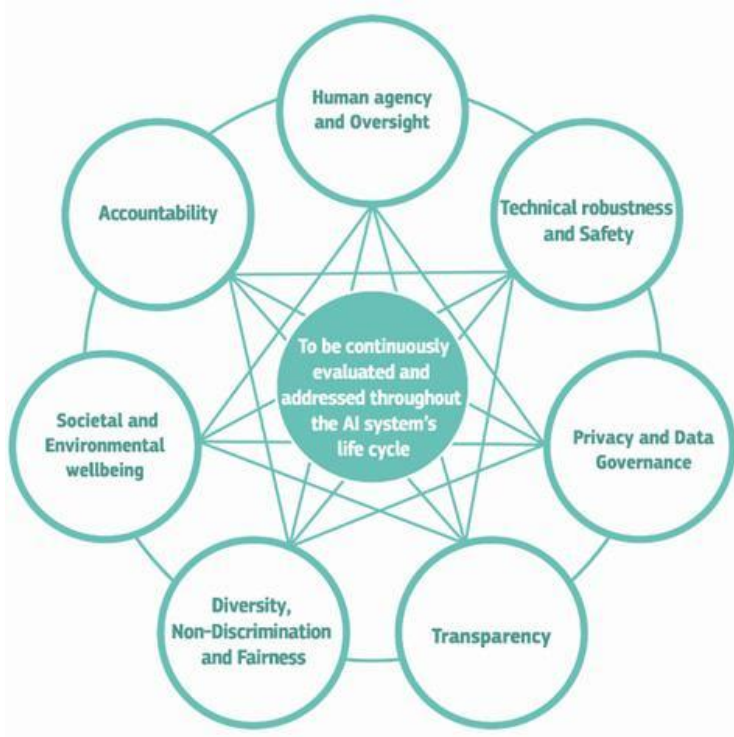
A curated hub of frameworks, standards, and educational tools

Welcome to the AI Governance Library, your self-service resource hub, bringing together hands-on materials to support your journey.

AI Governance is dynamic and evolving. Our resources enable you to revisit and refine your AI governance strategies, ensuring they remain effective, relevant, and compliant over time. We will keep adding resources as the technology and best practices evolve.



Question Bank - Considerations for Implementing AI



Fundamental rights encompass rights such as human dignity and non-discrimination, as well as rights in relation to data protection and privacy, to name just some examples. Prior to self-assessing an AI system with this Assessment List, a fundamental rights impact assessment (FRIA) should be performed.

A FRIA could include questions such as the following – drawing on specific articles in the Charter and the European Convention on Human Rights (ECHR)¹⁴ its protocols and the European Social Charter.¹⁵

1. Does the AI system potentially negatively discriminate against people on the basis of any of the following grounds (non-exhaustively): sex, race, colour, ethnic or social origin, genetic features, language, religion or belief, political or any other opinion, membership of a national minority, property, birth, disability, age or sexual orientation?

Have you put in place processes to test and monitor for potential negative discrimination (bias) during the development, deployment and use phases of the AI system?

Have you put in place processes to address and rectify for potential negative discrimination (bias) in the AI system?

2. Does the AI system respect the rights of the child, for example with respect to child protection and taking the child's best interests into account?

Have you put in place processes to address and rectify for potential harm to children by the AI system?

Have you put in place processes to test and monitor for potential harm to children during the development, deployment and use phases of the AI system?

Have an AI learners within your organisation

Their Job: ½ day per week (10% of their time)

- Reading news articles, Newsletters
- Go to meetups, conferences
- Take free online courses from DeepLearning AI

Suggested Output (Approx 3- 6 months)

**Write an AI Policy for
your organisation**

**Automate 1 process
that can save their team
10 days of work**

3 FEBRUARY 2026 — ANNUAL REPORT

International AI Safety Report 2026

The second International AI Safety Report, published in February 2026, is the next iteration of the comprehensive review of latest scientific research on the capabilities and risks of general-purpose AI systems. Led by Turing Award winner Yoshua Bengio and authored by over 100 AI experts, the report is backed by over 30 countries and international organisations. It represents the largest global collaboration on AI safety to date.

Translated versions in the other 5 official UN languages can be found under the 'More Languages' button. The 'Extended Summary for Policymakers' can be found on the main 'Publications' page.



INTERNATIONAL AI SAFETY REPORT

2026

**Extended Summary
for Policymakers**

February 2026



“Both optimists and pessimists contribute to society. The optimist invents the aeroplane, the pessimist the parachute.”

- George Bernard Shaw, playwright, critic, political activist

New Zealand's first AI Safety Conference 4th & 5th July



A screenshot of the event registration page for the AI Safety New Zealand Conference 2026. The page has a light blue background. At the top left is the 'AIS ANZ' logo. The main title is 'AI Safety New Zealand Conference 2026'. The event is scheduled for Saturday 4 July, from 9:30 - 5:30, 16:30, at EPIC Christchurch (Enterprise Precinct and Innovation Campus) in Christchurch, Canterbury Region. The page lists ticket options: 'Full Conference (Sat & Sun)' for NZ\$200.00, 'Day 1: Theory & Foundations' for NZ\$120.00, and 'Day 2: Action & Impact' for NZ\$120.00. There is also an 'Academic & Researcher Rate' for NZ\$100.00. A 'Request to Join' button is at the bottom right.

Register Here  luma.com/AI_Safety_NZ



“New Zealand has the potential to be the smartest adopter of high risk, integrated AI technologies and be the playbook for other countries on how to safely implement these technologies”

**- Emma Humphrey, Strategist | Social Entrepreneur
Consultant on AI Safety & Responsible AI Ecosystem in NZ**