



Pacific AI

# Pacific AI Governance Policy Suite:

June 2026 Release Notes

<b>Introduction .....</b>	<b>7</b>
1. Purpose .....	7
2. Scope .....	7
3. Glossary .....	8
4. How to Use this AI Policy Suite .....	8
5. Your License to this AI Policy Suite .....	10
6. Disclaimer .....	10
<b>Covered Laws, Regulations, Frameworks &amp; Standards .....</b>	<b>12</b>
1. Frameworks and Standards .....	12
2. Healthcare Guideline Frameworks .....	12
3. International Treaties .....	13
4. US National Legislation .....	13
5. US Federal Regulation .....	14
6. US State & Local Legislation .....	14
7. US State & Local Legislation – Privacy Laws .....	16
8. US State & Local Legislation – Deepfake Laws .....	16
9. US State & Local Legislation – Healthcare AI Laws .....	17
10. US State & Local Legislation –Finance and Insurance .....	18
11. US State & Local Legislation – Employment and Recruitment .....	19
12. The European Union .....	19
13. The United Kingdom .....	19
14. Argentina .....	20
15. Australia .....	20
16. Brazil .....	21

17.	Canada .....	21
18.	Colombia .....	21
19.	India .....	22
20.	Indonesia.....	22
21.	Israel.....	22
22.	Japan .....	22
23.	Mexico.....	23
24.	Norway.....	23
25.	Saudi Arabia .....	23
26.	Singapore .....	23
27.	South Korea.....	24
28.	Switzerland .....	24
29.	Taiwan.....	24
30.	Thailand .....	25
31.	Turkey .....	25
32.	UAE .....	25
33.	Acceptable Use Policies by Major Providers .....	26
34.	Contractual Clauses Checklists.....	26
	<b>AI Risk Management Policy.....</b>	<b>27</b>
1.	Purpose .....	27
2.	Scope.....	27
3.	AI Governance Officer .....	27
4.	Risk Management.....	27
5.	Risk Evaluation .....	29

6.	Risk Mitigation.....	30
7.	Risk Management Schedule .....	32
8.	Vendor Management .....	33
9.	Training and AI Literacy.....	34
10.	Reporting .....	35

## AI System Lifecycle Policy ..... 36

1.	Purpose .....	36
2.	Scope.....	36
3.	Risk Manager.....	36
4.	Go / No-Go Checkpoint .....	37
5.	Risk Level.....	38
6.	Pre-Deployment Checkpoint .....	39
7.	System Development and Testing.....	40
8.	Declaration of Conformity for High-Risk Systems .....	42
9.	Post Market Surveillance and Annual Checkpoint .....	43
10.	Termination or Deletion of models or AI systems .....	44
11.	Inventory for Internal and External AI Systems .....	44

## AI Safety Policy ..... 45

1.	Purpose .....	45
2.	Scope.....	45
3.	Safety Threat Modeling.....	45
4.	Automated Testing.....	46
5.	Manual Testing.....	47
6.	Monitoring .....	48

- 7. Human Oversight & Override ..... 48
- 8. Additional Safety Controls for GPAI models ..... 49

## AI Privacy Policy ..... 50

- 1. Purpose ..... 50
- 2. Scope ..... 50
- 3. Privacy by Design ..... 50
- 4. Profiling ..... 51
- 5. Informed, Specific, and Revocable Consent for Personal Data Use ..... 51
- 6. De-Identification of Training Data ..... 52
- 7. Testing for Privacy ..... 53
- 8. Child Privacy ..... 53
- 9. Incident Reporting ..... 53

## AI Fairness Policy ..... 55

- 1. Purpose ..... 55
- 2. Scope ..... 55
- 3. Bias Threat Modeling ..... 55
- 4. Mitigation of Data Bias ..... 57
- 5. Mitigation of Algorithmic Bias ..... 59
- 6. Testing for Data and Algorithmic Bias ..... 60
- 7. Manual Testing and Bias Audits ..... 61
- 8. Mitigating for Unfair Outcomes in Production ..... 62
- 9. Monitoring in Production ..... 63

## AI Transparency Policy ..... 65

- 1. Purpose ..... 65

2.	Scope .....	65
3.	Disclose Use of AI .....	65
4.	Disclose AI Generated Content .....	66
5.	No Deepfakes .....	66
6.	Explain AI Decisions.....	66
7.	Disclosures When Acting as an AI Developer .....	67
8.	Disclosures When Acting as an AI Deployer .....	69
<b>AI Incident Reporting Policy.....</b>		<b>70</b>
1.	Purpose .....	70
2.	Scope .....	70
3.	Internal Reporting .....	70
4.	External Reporting.....	71
5.	Classification of Incidents.....	72
6.	Whistleblower Policy.....	73
<b>AI Copyright Policy.....</b>		<b>75</b>
1.	Purpose .....	75
2.	Scope.....	75
3.	Web Crawling Controls.....	75
4.	Mitigating the Risk of Infringing Outputs .....	76
5.	Complaint Mechanism .....	76
<b>AI Acceptable Use Policy.....</b>		<b>77</b>
1.	Purpose .....	77
2.	Scope.....	77
3.	Unacceptable Uses of AI.....	77
4.	Enforcement and Review .....	80

# Introduction

## 1. Purpose

The goal of this suite of organizational policies is to establish a unified set of policies and guidelines for artificial intelligence (AI) systems. Subject to our Disclaimer Section below and together with continuous operational, technical and organizational measures, they help ensure your organization conforms to all relevant laws, regulations, and industry standards.

The policies define roles & responsibilities required to build and operate AI systems that are legal – as well as safe, effective, fair, transparent, accountable, private, and secure.

## 2. Scope

This Policy Suite aims to enable your organization to develop or deploy AI systems legally anywhere in the US, EU and certain other countries. This document lists the current list of sources which it covers. If you find a missing or outdated source, please email [legal@pacific.ai](mailto:legal@pacific.ai) with the details so that it can be included in the next version of this suite. Updates are released on a quarterly basis. The inclusion criteria into this policy suite are as follows:

1. Enacted legislation is included. The policies include provisions for addressing a requirement if it is legislated somewhere, so that by using these policies we can legally operate anywhere in the covered countries. This means that the policies may be more restrictive than legally necessary, if you're strictly operating in a subset of US states. However, since most AI systems are online and serve customers across state lines, unifying all the requirements across checkered legislation is required.
2. Enacted regulation is included. For example, guidelines and rules issued by US regulators such as the FTC, FDA, and EEOC are legally enforceable, and hence your organization must comply with them.
3. Published international standards are included. Standards from organizations like NIST and ISO are sometimes used in court cases to define 'commercially reasonable' efforts to manage or mitigate a risk. Therefore, by complying with these standards, you are making reasonable efforts to address these risks.
4. Other publications, which provide general guidance or best practices, are excluded. This includes for example the US AI Bill of Rights or the Bletchley Declaration. While these define national intentions, they are not legally binding.
5. Legislation that has not passed yet (but only proposed) is excluded. Similarly, standards and frameworks still under development are excluded. These will be included in future versions of these policies when they become law, regulation, or industry standard.

### 3. Glossary

An “**AI system**” means a machine-based system that is designed to operate with varying levels of autonomy and that may exhibit adaptiveness after deployment, and that, for explicit or implicit objectives, infers, from the input it receives, how to generate outputs such as predictions, content, recommendations, or decisions that can influence physical or virtual environments.

A “**GPAI model**” (General-Purpose AI model) means an AI model that is trained with a large amount of data using self-supervision at scale, that displays significant generality and is capable of competently performing a wide range of distinct tasks regardless of the way the model is placed on the market, and that can be integrated into a variety of downstream systems or applications.

An “**AI developer**” is a provider – whether a person, corporation, public authority, agency, or other legal entity that develops an AI system or general-purpose AI model and markets it under its own name or trademark. This includes those who develop AI or have it developed by others, and who then market it under their name.

An “**AI deployer**” is as any individual or entity that uses an AI system within their professional scope, excluding personal and non-professional activities. This expansive definition includes any business use of AI.

### 4. How to Use this AI Policy Suite

The suite contains nine policies that your organization needs to adopt and implement to conform to the covered legislation, regulations, and standards:

- AI Risk Management Policy
- AI System Lifecycle Policy
- AI Safety Policy
- AI Privacy Policy
- AI Fairness Policy
- AI Transparency Policy
- AI Incident Reporting Policy
- AI Copyright Policy
- AI Acceptable Use Policy

Most controls in this Policy Suite apply to both AI systems and GPAI models, whether you are acting as a developer or deployer. Sections that only apply in some cases are titled as such.

To put these policies to use:

- **Adopt the policies.** Read them, customize them as needed, and formally adopt them as required policies across your organization. This requires the written approval of your cross-functional team, executive team and the appointment of an AI Governance Officer.
- **Formal approval.** The AI Governance Policy Suite should be approved by organizational leadership.
- **Implement the policies.** This involves a range of controls from training, risk management, incident reporting, safety and bias threat modeling, monitoring of production systems, and other controls which the policies specify.
- **Attest that you comply.** Email [info@pacific.ai](mailto:info@pacific.ai) with a written confirmation that the policies are adopted, approved, and implemented, and we'll provide you with a website badge you can use to publicly show self-attestation for organization-wide AI Governance.

You can do this without having to pay anything to Pacific AI, as long as you abide by the license terms explained in the next section. If you do need help implementing the policies, or are required to have a third-party certification, Pacific AI offers these as paid services. Each project is customized to fit your organization's needs, size, and starting point.

Email [info@pacific.ai](mailto:info@pacific.ai) to start a discussion about:



#### AI Governance Implementation

Get expert help building enterprise AI governance aligned with NIST and CHAI.



#### Safe GenAI Deployment

Work with your team to test, evaluate, and monitor generative AI — safely and at scale.



#### AI Policy Suite

A free, ready-to-use AI policy framework aligned with over 100 global healthcare regulations and standards.



#### AI Audit

A third-party professional audit service to assess your AI system's accuracy, safety, bias, and regulatory compliance.

## 5. Your License to this AI Policy Suite

The AI Policy Suite is provided under the [CC BY-NC-SA 4.0](#) License.

Under this license, you are free to:

- Use the content at no cost within your organization.
- Share, copy and redistribute the material in any medium or format.
- Adapt, remix, transform, customize, and build upon the material.

Pacific AI cannot revoke these freedoms if you follow the license terms, which include:

- Attribution. You must give appropriate credit, retain the logo of Pacific AI, provide a link to Pacific AI and to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests Pacific AI endorses you or your use.
- No Commercial Use. Selling, reselling, publishing to attract traffic, or embedding the content in a product or service are not. Internal use is allowed, such as sharing the content within the organization or with customers, partners, or investors.
- Share-Alike. If you adapt or build upon the material, you must distribute your contributions under the same license.
- No additional restrictions. You may not apply legal terms or technological measures that legally restrict others from doing anything the license permits.

## 6. Disclaimer

The suggestions found in this Policy are provided for information purposes to guide the responsible development, deployment, and use of AI systems within your organization. While every effort has been made to ensure that the principles, procedures, and controls described herein reflect best practices and applicable regulatory requirements at the time of publication, there is no guarantee of full compliance with all applicable laws, regulations, or industry standards. Laws and regulatory frameworks governing AI are evolving rapidly, and it is the responsibility of each business unit to consult qualified legal counsel and compliance experts as needed.

AI technologies, their applications, and associated risks change over time. This Policy Suite may become outdated as technology, or legal requirements evolve. Interim changes in applicable law, ethical norms, or technical standards may require immediate updates not yet reflected herein.

Adoption of this Policy Suite alone does not constitute compliance with any applicable law, regulation, or industry standard. Compliance requires a company to implement, maintain, and continuously monitor the operational, technical, and organizational measures described herein.

Nothing in this Policy Suite, nor in any AI system, tool, or output referenced here, is intended to provide, or should be construed as providing:

- Legal advice or interpretation of laws, regulations, contracts, or other binding obligations;
- Medical advice, diagnoses, treatment recommendations, or any form of health-related guidance;
- Any other form of regulated professional advice requiring a license, certification, or statutory authority.

All decisions in these and other regulated areas must be made in consultation with appropriately qualified and authorized professionals.

The company remains responsible for taking all necessary actions, exercising appropriate diligence, and ensuring that its practices, systems, and personnel adhere to all legal and regulatory obligations. Nothing in this Policy Suite should be construed as a guarantee of compliance or as a substitute for ongoing risk management, legal review, or regulatory oversight.

Pacific AI shall not be liable for any loss, damage, or claim arising directly or indirectly from reliance on this policy suite. In the event of any conflict between this Policy Suite and applicable laws or regulations, the latter shall prevail.

# Covered Laws, Regulations, Frameworks, and Standards

## 1. Frameworks and Standards

- NIST AI Risk Management Framework
- NIST AI Risk Management Framework – Generative AI Profile
- ISO/IEC 42001:2023 - AI management systems
- ISO/IEC 42005:2025 - AI system impact assessment
- ISO/IEC 23894:2023 - Guidance on AI risk management
- ISO/IEC NP TS 12831 - Testing for AI Systems
- ISO/IEC 5259-4:2024 - Data quality for analytics and machine learning
- NIST Privacy Framework (PF)
- NIST Cybersecurity Framework (CSF)
- American Bar Association (ABA) Guidance on AI in Employment Law
- IEEE Standard for the Procurement of Artificial Intelligence and Automated Decision Systems

## 2. Healthcare Guideline Frameworks

- Coalition for Health AI (CHAI) – Assurance Standards Guide & AI Governance Playbooks
- FUTURE-AI - Consortium international consensus guideline for trustworthy and deployable artificial intelligence in healthcare
- QUEST - Framework for Human Evaluation of LLMs in Healthcare Derived from Literature Review
- CLAIM - Checklist for Artificial Intelligence in Medical Imaging
- TEHAI - Translational Evaluation of Healthcare AI
- MEDIC - Comprehensive Framework for Evaluating LLMs in Clinical Applications
- STARD-AI - Standards for Reporting of Diagnostic Accuracy Study
- MI-CLAIM-GEN - Minimum Information about Clinical Artificial Intelligence Checklist for Generative Modelling Research
- CRAFT-MD - Conversational Reasoning Assessment Framework for Testing in Medicine
- TRIPOD-LLM - Transparent reporting of a multivariable model for individual prognosis or diagnosis

- TRIPOD-AI - Updated guidance for reporting clinical prediction models that use regression or machine learning methods
- CONSORT-AI - Reporting guidelines for clinical trial reports for interventions involving AI
- SPIRIT-AI – Guidelines for clinical trial protocols for interventions involving AI
- AMA - Augmented Intelligence Development, Deployment, and Use in Health Care
- WHO – Ethics and governance of artificial intelligence for health
- WHO – Generating Evidence for Artificial Intelligence Based Medical Devices: A Framework for Training Validation and Evaluation
- HAIRA – Advancing Healthcare AI Governance: A Comprehensive Maturity Model Based on Systematic Review
- TPLC – Total Product Lifecycle framework for Healthcare AI/ML
- OPTICA – Organizational Perspective Checklist for AI solutions adoption
- SALIENT – Implementation frameworks for end-to-end clinical AI
- AHRQ & AIMHD Guiding principles to address the impact of algorithm bias
- ‘Model Facts’ label for HTI-1 compliance by Duke Institute for Health Innovation
- IMDRF – “Software as a Medical Device”: Possible Framework for Risk Categorization and Corresponding Considerations
- National Academy of Medicine – Health Care Artificial Intelligence Code of Conduct

### 3. International Treaties

- OECD Framework for the Classification of AI systems
- OECD Common Reporting Framework for AI Incidents
- CSET’s AI Incidents Key Components for a Mandatory Reporting Regime
- UNESCO’s Ethical Impact Assessment
- UNESCO Recommendation on Ethics of AI
- Council of Europe Framework Convention on Artificial Intelligence
- UNICEF’s Guidance on AI and Children
- Global Digital Compact

### 4. US National Legislation

- Title VII of the Civil Rights Act of 1964
- Section 1981 of the Civil Rights Act of 1866

- The Equal Pay Act of 1963
- The Age Discrimination in Employment Act of 1967
- The Immigration Reform and Control Act
- Titles I and V of the Americans with Disabilities Act of 1990 (ADA)
- The Pregnant Workers Fairness Act
- The Uniformed Services Employment and Reemployment Rights Act
- The Genetic Information Nondiscrimination Act
- The Take it Down Act (S.146)

## 5. US Federal Regulation

- [ACA 1557](#): Nondiscrimination in Health Programs and Activities
- [HHS HTI-1](#): Health Data, Technology, and Interoperability
- [EEOC Guidance](#) on the Use of AI to Assess Job Applicants and Employees
- [CFPB Guidance](#) on Black-Box Credit Models
- [FDA Draft Guidance on Considerations for the Use of AI to Support Regulatory Decision-Making for Drug and Biological Products](#)
- FDA Draft Guidance for Developers of AI-Enabled Medical Devices
- FDA Final Rule codifying 21 CFR Part 892 — “Classification of the Radiological Computer-Assisted Detection and Diagnosis Software”
- Executive Order on Removing Barriers to American Leadership in AI
- America’s AI Action Plan
- AI Literacy Framework
- OMB Memorandum M-25-22 on Driving Efficient Acquisition of AI in Government
- Executive Order on Ensuring a National Policy Framework for Artificial Intelligence
- Executive Order on Advancing Artificial Intelligence Education for American Youth
- Executive Order on Promoting Advanced Artificial Intelligence Innovation and Security

## 6. US State & Local Legislation

- [Arkansas Act 927](#): Establishing Ownership Rights for Content Generated and Models Trained Using Generative AI Tools
- California SB53: Transparency in Frontier Artificial Intelligence Act

- California SB243: Companion Chatbots
- California [AB2013: generative AI training data transparency](#)
- California [AB2855: Artificial intelligence Law](#)
- California [SB-896: Generative Artificial Intelligence Accountability Act](#)
- California SB-942 and AB-853: AI Transparency Act
- California SB-1120 Health Care Coverage: Utilization Review
- California AB2885: Unified Definition of Artificial Intelligence
- California SB361: Data Broker Disclosure Requirements to CPPA
- California AB316: AI defenses in litigation
- Colorado AI Act [to be replaced by Colorado SB 26-189]
- Colorado SB26-189: Automated Decision-Making Technology
- Colorado HB 1263: Conversational Artificial Intelligence Service Operator Requirements
- Connecticut SB 5: AI Responsibility and Transparency Act
- Georgia SB 540: Online Internet Safety; Certain Disclosures Related to Conversational AI Services
- Idaho SB 1297: Conversational AI Safety Act
- [Illinois AI Video Interview Act](#)
- Illinois HB 3773: Illinois Human Rights Act
- Maine HP 1154: Ensure Transparency in Consumer Transactions Involving Artificial Intelligence
- Montana HB 178: Limit government use of AI systems
- Nebraska LB 525: Conversational AI Safety Act
- New York Comprehensive Healthcare Appeals Reform Act
- New York SB 7543: Legislative oversight of automated decision-making in government act
- New York RAISE Act
- New York SB 8420A: requiring advertisements to disclose the use of a synthetic performer
- [Illinois AI Video Interview Act](#)
- Illinois HB 3773: Illinois Human Rights Act
- Illinois SB 315: AI Safety Measures Act
- Maine HP 1154: Ensure Transparency in Consumer Transactions Involving Artificial

## Intelligence

- Montana HB 178: Limit government use of AI systems
- New York SB 7543: Legislative oversight of automated decision-making in government act
- Oregon SB 1546: relating to AI companions
- Texas HB 149: Responsible Artificial Intelligence Governance Act
- Texas SB1964: regulation and use of artificial intelligence systems and the management of data by governmental entities
- Washington HB 2225 and SB 1546: Chatbot Disclosure Act

## 7. US State & Local Legislation – Privacy Laws

- [California Consumer Privacy Act](#)
- [California Privacy Rights Act](#)
- [Colorado Privacy Act](#)
- Colorado [SB22-113](#): Artificial Intelligence Facial Recognition
- [Connecticut Data Privacy Act](#)
- [Delaware Personal Data Privacy Act](#)
- Indiana [SB5](#): Consumer data protection
- Minnesota [HF2309](#): Consumer rights provided regarding personal data
- Minnesota [HF 4757](#): Minnesota Consumer Data Privacy Act
- Montana [Consumer Data Privacy Act](#)
- New Hampshire [SB 255](#): Expectation of privacy
- Oregon [SB619](#): Protections for the personal data of consumers
- [Texas Data Privacy and Security Act](#)
- [Utah Consumer Privacy Act](#)
- [Utah AI Policy Act](#)
- [Virginia Consumer Data Privacy Act](#)

## 8. US State & Local Legislation – Deepfake Laws

- Arizona HB2394: Digital Impersonation
- Arizona SB1359: Prohibition of Deepfakes on Election Communications
- Arkansas Act 977: Concerning Criminal Offenses Related to Possession of Sexually Explicit

#### Material that Depicts a Child

- [California Bolstering Online Transparency Act](#)
- California [AB2602](#): Contracts against public policy
- California [AB1836](#): Use of likeness: digital replica
- California [AB621](#): Deepfake pornography
- Colorado [HB1147](#): Candidate Election Deepfake Disclosures
- Florida [HB 919](#): Artificial Intelligence Use in Political Advertising
- Florida [SB 1798](#): Sexually Related Offenses
- [Georgia § 16-11-90](#): Prohibition on Nude or Sexually Explicit Electronic Transmissions
- [Hawaii SB 309](#): Privacy in the First Degree
- [Illinois HB 2123](#): Nonconsensual Dissemination of Private Sexual Images Act
- [Illinois HB 3773](#): Limit Predictive Analytics Use
- [Minnesota HB 1370](#): nonconsensual dissemination of deep fake sexual images
- Nebraska [LB383](#): Prohibition of Generated Child Sexual Abuse Material in Nebraska statutes
- New Jersey [A3540](#) and [S2544](#): Establishing Civil and Criminal Penalties for Production and Dissemination of Deceptive Audio and Visual Media
- [New York S1042A](#): unlawful dissemination or publication of intimate images
- [South Dakota SB120](#): Record, Privacy, Manipulated image violation
- Tennessee [ELVIS Act](#)
- [Texas PENAL § 21.165](#): Unlawful Production or Distribution of Certain Videos
- [Texas SB441](#): Unlawful creation of sexually explicit deepfakes
- [Texas HB 581](#): Age verification requirement or prohibition of AI-generated sexual material harmful to minors
- [Virginia § 18.2-386.2](#): Unlawful dissemination or sale of images of another
- Washington [HB1205](#): Prohibiting the Knowing Distribution of a Forged Digital Likeness

## 9. US State & Local Legislation – Healthcare AI Laws

- Arizona [HB 2175](#): Prohibition of the use of AI to deny claims involving medical judgment
- California [SB 1120](#): Regulation of health plans using AI for utilization review
- California [AB 3030](#): Requirement of AI-generated patient communications to include a disclaimer and instructions to contact a human provider

- California AB 489: Prohibiting AI systems from Misrepresenting Themselves as Licensed Healthcare Professionals
- Colorado HB 1139: Use of AI in Healthcare
- Colorado HB 1195: Psychotherapy AI Restrictions
- Illinois HB1806: Wellness and Oversight for Psychological Resources Act
- Maryland HB 820: Requiring carriers to ensure that any AI tools used for utilization review base decisions on medical / clinical history, individual circumstances, and clinical information
- Nebraska LB 77: Ensuring Transparency in Prior Authorization Act and provide for insurance and Medicaid coverage of biomarker testing
- Nevada AB 406: Prohibition of AI providers from indicating that an AI system can provide professional mental or behavioral health care
- Oregon HB 2748: Prohibition of non-human entities to use the title of nurse or similar titles
- Tennessee SB 1580: Health care AI act
- Texas SB1188: Prohibition of the offshoring of electronic medical records, and implement safeguards for patient data
- Texas SB 815: use of certain automated systems in, and certain adverse determinations made with the health benefit claims process
- Utah HB 452: Suppliers of mental health chatbots are required to disclose that the chatbot is AI technology
- Virginia HB 2154: Requirement of hospitals, nursing homes, and certified nursing facilities to establish and implement policies on access to, and use of, an intelligent personal assistant at their facility

## 10. US State & Local Legislation –Finance and Insurance

- Equal Credit Opportunity Act
- Fair Credit Reporting Act
- Federal Trade Commission Act
- California AB325: Cartwright Act: violations
- Colorado [SB 21-169](#) Restrict Insurers' Use of External Consumer Data
- Indiana HB 1271: Payment of Health Claims
- Utah SB 319: Health Insurance Preauthorization Amendment
- Washington SB 5395: Transparency and accountability in the prior authorization process

- New York State Department of Financial Services Insurance Circular Letter No 7 on Use of AI Systems and External Consumer Data and Information Sources in Insurance Underwriting and Pricing
- SEC AI and Investment Fraud: Investor Alert
- NAIC Model Bulletin: Use of AI Systems by Insurers (adopted in approx. 30 states)

## 11. US State & Local Legislation – Employment and Recruitment

- California SB7: Employment; Automated Decision Systems
- New York SB822: Disclosure of automated employment decision-making tools and maintaining an artificial intelligence inventory
- Final Employment Regulations Regarding Automated Decision Systems
- New York City Local Law 144: Automated employment decision tools

## 12. The European Union

- The EU AI Act
- The General-Purpose Code of Practice
- GDPR
- Digital Services Act
- Digital Markets Act
- Product Liability Directive
- Ethics guidelines for trustworthy AI
- Guidelines on AI system definition
- Guidelines on prohibited AI practices
- Guidelines for providers of general-purpose AI models
- Guidelines on the classification of high-risk AI systems

## 13. The United Kingdom

- AI Regulation White Paper by the UK Government (March 2023)
- A pro-innovation approach to AI regulation: government response (February 2024)
- Data (Use and Access) Act 2025
- Data Protection Act 2018
- Human Rights Act 1998

- Equality Act 2010
- Online Safety Act 2023
- Medical Devices Regulations 2002
- Advertising Standards Authority: Generative AI & Advertising: Decoding AI Regulation
- CMA: AI Foundation Models: Initial Report
- CQC: Using machine learning in diagnostic service: A report with recommendations from CQC's regulatory sandbox
- EHRC: EHRC guidance on use of AI by public bodies
- ICO: Guidance on AI and data protection
- ICO: Explaining decisions made with AI
- MHRA: AI Roadmap and AI Airlock regulatory sandbox
- FCA: Digital and AI sandboxes
- NCSC: AI cybersecurity guidelines
- Ofcom: Synthetic media (including deepfakes) in broadcast programming
- OPSS: Study on the impact of AI on product safety
- Office for Artificial Intelligence: Guidelines for AI procurement

## 14. Argentina

- Data Protection Law
- Recommendations for Reliable AI
- Program for Transparency and Data Protection in the Use of AI
- Responsible AI Guide

## 15. Australia

- AI Ethics Principles
- Guidance for AI Adoption
- Online Safety Act
- AI Policy Guide and Template
- National AI Plan
- Model AI Governance Framework for Agentic AI
- Guidance for AI Adoption

- Policy for the Responsible Use of AI in Government
- Standard for AI Transparency Statements for Australian Government Agencies

## 16. Brazil

- General Data Protection Law
- Consumer Protection Code
- Copyright Law
- Law 15.123/2025 on Aggravating Penalty for Crimes Using AI Deepfakes against Women
- Federal Public Sector Ethical AI Impact Assessment Framework
- Federal Government Generative AI Guidance
- National Council of Justice AI Rules

## 17. Canada

- Privacy Act and the Personal Information Protection and Electronic Documents Act (PIPEDA)
- Canadian Human Rights Act
- Directive on Automated Decision-Making
- Pan-Canadian Artificial Intelligence Strategy
- Alberta Personal Information Protection Act, SA 2003
- Alberta Protection of Privacy Act
- British Columbia Personal Information Protection Act
- Quebec Act Representing the Protection of Personal Information in the Private Sector
- Quebec provincial privacy law with explicit automated-decision obligations

## 18. Colombia

- National Artificial Intelligence Policy
- General Data Protection Framework
- Consumer Protection Statute
- National AI Policy 2024-2030
- Ethical Guide for AI in Public Entities
- Information Security and Privacy Guidelines for AI Systems

## 19. India

- AI Governance Guidelines
- National Strategy for Artificial Intelligence
- Principles for Responsible AI
- Operationalizing Principles for Responsible AI
- Information Technology Act
- Digital Personal Data Protection Act
- IT (Intermediary Guidelines and Digital Media Ethics Code) Amendment Rule
- IT Rules Amendment on AI-Generated Content and Deepfakes

## 20. Indonesia

- National AI Strategy
- Circular Letter No. 9/2023 on Artificial Intelligence Ethics
- Responsible and Trustworthy AI Code of Conduct for fintech
- Personal Data Protection Law, PDP Law

## 21. Israel

- Policy on AI Regulation and Ethics
- Copyright Act
- Protection of Privacy Law
- Privacy-Enhanced Technologies Guidelines

## 22. Japan

- Act on Promotion of Research and Development and Utilization of Artificial Intelligence-Related Technologies
- Digital Platform Transparency Act
- Financial Instruments and Exchange Act
- Copyright Act
- Act on the Protection of Personal Information
- Information Distribution Providers Act
- Social Principles of Human-Centric AI

- AI Guidelines for Business by METI and MIC
- AI Institutional Council Guidelines

## 23. Mexico

- Federal Law on Protection of Personal Data Held by Individuals
- National AI Agenda for 2024–2030
- National Declaration on Ethical AI

## 24. Norway

- Equality and Discrimination Act
- Working Environment Act
- Transparency Act
- Personal Data Act
- Marketing Control Act
- Consumer Purchase Act
- Digital Services Act
- Position Paper on the European Commission’s Proposal for a Regulation of the European Parliament and of the Council Laying Down Harmonized Rules on AI
- Norway’s National AI Strategy

## 25. Saudi Arabia

- AI Ethics Principles
- AI Readiness Index
- Personal Data Protection Law
- Generative AI Guidelines
- General Rules for Secondary Use of Data

## 26. Singapore

- Model AI Governance Framework
- Model AI Governance Framework for Generative AI
- Model AI Governance Framework for Agentic AI
- AI Verify Toolkit

- National Artificial Intelligence Strategy 2.0
- Health Products Act 2007
- Advisory Guidelines on Use of Personal Data in AI Recommendation and Decision Systems
- Principles to Promote Fairness, Ethics, Accountability and Transparency (FEAT) in the Use of Artificial Intelligence and Data Analytics in Singapore’s Financial Sector
- Artificial Intelligence in Healthcare Guidelines

## 27. South Korea

- South Korea AI Basic Act (Basic Act on the Development of Artificial Intelligence and the Establishment of Foundation for Trustworthiness)
- Act on Promotion of Information and Communications Network Utilization and Information Protection
- Personal Information Protection Act
- Fair Hiring Procedure Act
- Content Industry Promotion Act
- Copyright Act
- Public Official Election Act

## 28. Switzerland

- Federal Act on Data Protection
- Federal Copyright Act
- Gender Equality Act
- Disability Discrimination Act
- Federal Criminal Code
- Product Liability Act
- Guidelines on Artificial Intelligence for the Confederation
- FINMA Guidance on Governance and risk management when using artificial intelligence
- Digital Switzerland Strategy 2025

## 29. Taiwan

- AI Basic Act
- AI Technology R&D Guidelines

- Taiwan Artificial Intelligence Action Plan

## 30. Thailand

- Personal Data Protection Act
- Computer Crime Act
- AI Ethics and Principles Guidelines
- Generative AI Governance Guidelines for Organizations
- National AI Strategy
- Securities and Exchange Commission Framework for AI and ML in Capital Markets

## 31. Turkey

- National AI Strategy
- Recommendations on the Protection of Personal Data in the Field of Artificial Intelligence
- Guidelines on Protection of Privacy in Mobile Applications
- Guidelines on Generative AI

## 32. UAE

- UAE AI Act
- National Strategy / Artificial Intelligence 2031 Initiative
- Charter for the Development and Use of Artificial Intelligence
- Personal Data Protection Law 2021
- Project of Future Nature (Law 25 of 2018)
- Deepfake Guide (2021)
- AI Ethics Guide (2022)
- AI Adoption Guideline in Government Services (2023)
- Responsible Metaverse Self-Governance Framework (2023)
- Guidelines for Financial Institutions adopting Enabling Technologies
- Federal Decree-Law No. 38 of 2021 on Copyright and Neighboring Rights (the UAE Copyright Law)
- DIFC Data Protection Regulations
- Abu Dhabi Department of Health's Policy on the Use of Artificial Intelligence in the Healthcare Sector

- Abu Dhabi Department of Health's Responsible Artificial Intelligence Standard
- Central Bank of the UAE's Guidelines for Financial Institutions adopting Enabling Technologies

### 33. Acceptable Use Policies by Major Providers

- OpenAI Usage Policy
- Anthropic Usage Policy
- Microsoft Enterprise AI Services Code of Conduct
- AWS Responsible AI Policy
- Google Generative AI Prohibited Use Policy
- Meta Seamless Acceptable Use Policy
- Cohere Labs Acceptable Use Policy
- John Snow Labs AI Acceptable Use Policy
- Mistral Usage Policy
- Snowflake Acceptable Use Policy
- Databricks Open Model Acceptable Use Policy

### 34. Contractual Clauses Checklists

- EU AI Model Contractual Clauses by EU Community of Practice on Public Procurement of AI
- Australian Government AI Model Clauses
- Government of Japan Checklist for Contracts Concerning Development and Use of AI
- Vischer Checklist on AI for contracts with suppliers and partners
- Society for Computers and Law EU AI Act Contractual Clauses
- New Zealand AI Procurement Checklist
- FS-ISAC Generative AI Vendor Risk Assessment Guide
- UK Government Guidelines for AI Procurement

# AI Risk Management Policy

## 1. Purpose

This policy describes AI Risk Management processes at our organization (“we”, “us”, or “our”), and how we identify, manage, avoid, and report risks of AI systems, such as threats to human safety, dignity, or liberties.

## 2. Scope

This policy, in addition to applicable laws and regulations, is mandatory to us and our team members.

We expect our business partners to meet the same standards and obligations as we demand from ourselves. This policy applies to all team members (whether temporary, fixed-term, or permanent) and third parties (contractors, seconded staff, trainees) acting under our responsibility. Together referred to as "team members". The policy also applies to our officers, trustees, and/or management at any level.

Any arrangements our company makes with a third party are subject to clear contractual terms, including specific provisions that require the third party to comply with minimum standards and procedures set out in this policy or other policies.

## 3. AI Governance Officer

The role that is responsible for implementing this policy is the AI Governance Officer, who is appointed by the company’s chief executive officer, and has the minimum required qualifications for this role. The AI Governance Officer may delegate or assign tasks to others within the organization, establishing an “AI Governance Office”, while remaining responsible for their thorough completion. The AI Governance Officer has the primary responsibility for:

- Monitoring applicable laws and regulations
- Developing related training and courses
- Ensuring all covered team members comply with this policy
- Investigating allegations of improper activities
- Taking appropriate corrective and disciplinary actions
- Implementing an AI system impact assessment process

## 4. Risk Management

1. AI Risk analysis and risk management are recognized as important components of our corporate compliance program and information security program.

- Risk assessments are done throughout AI system life cycles;
  - Before the integration of new AI system and before changes are made to our physical safeguards;
  - These changes do not include routine updates to existing systems, deployments of new systems created based on previously configured systems, deployments of new customers, or new code developed for operations and management of our software;
  - We perform periodic technical and non-technical assessments of the security rule requirements.
2. Risk evaluation principles:
    - AI systems are tracked and monitored against negative risks, benefits and opportunities;
    - Use several lines of defense, where each team has different roles and independent responsibilities on each AI lifecycle level;
    - Encourage critical thinking and questioning throughout AI system life cycles;
    - Evaluate trade-offs using use cases and impacts.
  3. Security measures are implemented to reduce risks and vulnerabilities to a reasonable and appropriate level to:
    - Ensure the confidentiality, integrity, and availability of all AI systems that we use, maintain, or operate for our customers;
    - Protect against any reasonably anticipated threats or hazards to the security or integrity of customer data and infrastructure;
    - Ensure compliance by all team members.
  4. Any remaining (residual) risk after other risk controls have been applied, requires sign off by the AI Governance Officer.
  5. All our team members are expected to fully cooperate with all persons charged with doing risk management work, including contractors and audit personnel. Any team member that violates this policy will be subject to disciplinary action based on the severity of the violation.
  6. The implementation, execution, and maintenance of the information security risk analysis and risk management process is the responsibility of the AI Governance Officer (or other employee designated by them), and the identified Risk Management Team.
  7. All risk management efforts, including decisions made on what controls to put in place as well as those to not put into place, are documented and the documentation is maintained for six years.

8. The details of the AI Risk Management Process, including risk assessment, discovery, and mitigation, are outlined in detail below. The process is tracked, measured, and monitored using the following procedures:
  - The AI Governance Officer initiates the AI Risk Management Procedures by creating an Issue in our Incident Management System.
  - AI Compliance Officer is assigned to carry out the AI Risk Management Procedures.
  - Each finding is documented in an online database that is linked to each issue.
  - Once the AI Risk Management Procedures are complete, along with corresponding documentation, the AI Governance Officer approves or rejects the issue. If the issue is rejected, it goes back for further review and documentation.
  - If the review is approved, the AI Governance Officer then marks the issue as Done, adding any pertinent notes required.
9. The AI Risk Management Procedure is monitored at least on a quarterly basis using the Quality Management System reporting to assess compliance with the above policy.

## 5. Risk Evaluation

1. At least once per year, or as needed as part of incident response or risk re-evaluation, the AI Governance Officer will conduct an updated Risk Evaluation of AI related risks. This should reflect new research, case study, legislation, incident, feedback and follow-up– and its potential impacts on the organization’s risk profile.
2. The Risk Evaluation should engage multiple stakeholders, performed by an interdisciplinary team that reflect a wide range of capabilities, competencies, demographic groups, domain expertise, educational backgrounds, lived experiences, professions, and skills across the organization. It should also include people with different professional backgrounds: data scientists, experts, ethicists, human-computer interaction specialists, customers, software engineers, domain experts, users and user advocates, product managers, legal professionals, executives, and others.
3. The Risk Evaluation identifies potential benefits and harms, including less obvious or long-term impacts.
4. Potential harms related to AI systems that must be considered include:
  - harm to people, including harm to a person’s civil liberties, rights, physical or psychological safety, or economic opportunity; harm to a group such as discrimination against a population sub-group; harm to democratic participation or educational access;
  - harm to an organization, including harm to an organization’s business operations; harm

to an organization from security breaches or monetary loss; harm to an organization's reputation;

- harm to an ecosystem, including harm to interconnected and interdependent elements and resources;
- chemical, biological, radiological, or nuclear threats;
- harm to the global financial system, supply chain, or interrelated systems;
- cyber security threats;
- harm to natural resources, the environment, and the planet; broad potential AI harms, including immature safety or risk cultures related to Generative AI, public information integrity risks, impacts on democratic processes, or unknown long-term performance characteristics.

## 6. Risk Mitigation

Determination of appropriate controls to reduce risk is dependent upon the risk tolerance of the organization consistent with its goals and mission.

### Step 1. Prioritize Actions

- Using results from the AI Risk Evaluation, sort the threat and vulnerability pairs according to their risk-levels in descending order. This establishes a prioritized list of actions needing to be taken, with the pairs at the top of the list getting/requiring the most immediate attention and top priority in allocating resources
- Output - Actions ranked from high to low

### Step 2. Evaluate Recommended Control Options

- Although possible controls for each threat and vulnerability pair are arrived at in Step 7 of the AI Risk Assessment, review the recommended control(s) and alternative solutions for reasonableness and appropriateness. The feasibility (e.g., compatibility, user acceptance, etc.) and effectiveness (e.g., degree of protection and level of risk mitigation) of the recommended controls should be analyzed. In the end, select a "most appropriate" control option for each threat and vulnerability pair.
- Output - list of feasible controls

### Step 3. Conduct Cost-Benefit Analysis

- Determine the extent to which a control is cost-effective. Compare the benefit (e.g., risk reduction) of applying a control with its subsequent cost of application. Controls that are not cost-effective are also identified during this step. Analyzing each control

or set of controls in this manner, and prioritizing across all controls being considered, can greatly aid in the decision-making process.

- Output - Documented cost-benefit analysis of either implementing or not implementing each specific control

#### **Step 4. Select Control(s)**

- Taking into account the information and results from previous steps, the Risk Management Team determines the best control(s) for reducing risks to the information systems. These controls may consist of a mix of administrative, physical, and/or technical safeguards.
- Output - Selected control(s)

#### **Step 5. Assign Responsibility**

- Identify the workforce members with the skills necessary to implement each of the specific controls outlined in the previous step and assign their responsibilities. Also identify the training, and other resources needed for the successful implementation of controls. Resources may include time, money, equipment, etc.
- Output - List of resources, responsible persons and their assignments

#### **Step 6. Develop Safeguard Implementation Plan**

- Develop an overall implementation or action plan and individual project plans needed to implement the safeguards and controls identified. The Implementation Plan should contain the following information:
  - Each risk or vulnerability/threat pair and risk level;
  - Prioritized actions;
  - The recommended feasible control(s) for each identified risk;
  - Required resources for implementation of selected controls;
  - Team member responsible for implementation of each control;
  - Start date for implementation;
  - Target date for completion of implementation;
  - Maintenance requirements.
- The overall implementation plan provides a broad overview of the safeguard implementation, identifying important milestones and timeframes, resource requirements (staff and other individuals' time, budget, etc.), interrelationships

between projects, and any other relevant information. Regular status reporting of the plan, along with key metrics and success indicators should be reported to senior management.

- Individual project plans for safeguard implementation may be developed and contain detailed steps that assigned resources carry out to meet implementation timeframes and expectations. Additionally, consider including items in individual project plans such as a project scope, a list deliverable, key assumptions, objectives, task completion dates and project requirements.
- Output - Safeguard Implementation Plan

### Step 7. Implement Selected Controls

- As controls are implemented, monitor the affected system(s) to verify that the implemented controls continue to meet expectations. Elimination of all risk is not practical. Depending on individual situations, implemented controls may lower a risk level but not completely eliminate the risk.
- Continually and consistently communicate expectations to all Risk Management Team members, as well as senior management and other key people throughout the risk mitigation process. Identify when new risks are identified and when controls lower or offset risk rather than eliminate it.
- Additional monitoring is especially crucial during times of major environmental changes, organizational or process changes, or major facilities changes.
- If risk reduction expectations are not met, then repeat all or a part of the risk management process so that additional controls needed to lower risk to an acceptable level can be identified.
- Output - Residual Risk documentation

## 7. Risk Management Schedule

The two principal components of the AI risk management process - AI risk assessment and AI risk mitigation - will be carried out according to the following schedule to ensure the continued adequacy and continuous improvement of our AI governance program:

- On a Scheduled Basis - an overall risk assessment of our AI systems and governance will be conducted annually. The assessment process should be completed in a timely fashion so that risk mitigation strategies can be determined and included in the corporate budgeting process.
- Throughout an AI system's development lifecycle - from the time that a need for a new, untested information system configuration and/or application is identified through the time it is disposed of, ongoing assessments of the potential threats to a system and its vulnerabilities

should be undertaken as a part of the maintenance of the system.

- As Needed – the AI Compliance Officer may call for a full or partial risk assessment in response to changes in business strategies, information technology, information sensitivity, threats, legal liabilities, because of adversarial testing, or other significant factors that affect our organization.

## 8. Vendor Management

AI risk management applies both to systems built in-house as well as systems or components procured from third-party vendors. The policies and controls are the same whether an AI system being deployed was built in-house, bought or rented from a third party, or is a hybrid. Therefore:

1. Our legal and procurement teams must ensure that third-party vendors which supply AI systems or key components in such systems (such as models, services, datasets, or APIs) are contractually required to abide by equivalent AI Risk Management processes and procedures.
2. Vendor contracts must be reviewed prior to execution and on a regular basis to avoid arbitrary or capricious termination of critical AI technologies or vendor services and non-standard terms that may:
  - Amplify or defer liability in unexpected ways;
  - Enable or contribute to unauthorized data collection by vendors or third parties (e.g., secondary data use);
  - Give vendors or other third parties derivative rights in models, data, or other intellectual property.
3. Contracting with vendors should include requests for clauses regarding:
  - Warranties and meeting intended use;
  - Risk management and risk assessment;
  - Data and data governance, including obligations and restrictions on the use of data: no sale of data to third parties, no use of data for training vendor’s own AI models and other secondary use, no prohibited use of AI;
  - Notification and disclosure for serious incidents arising from third-party data and systems;
  - Service Level Agreements (SLAs) in vendor contracts that address incident response, response times, and availability of critical support;
  - Assignment of liability and responsibility for maintaining and updating security, privacy,

bias, component provenance, integrity checks, or safety incidents;

- Responsibility for monitoring and updating the AI system changes over time (e.g., fine-tuning, drift, decay, feedback);
- Responsibility for complying with AI laws, especially with respect to high-risk systems;
- Transparency and indemnification regarding intellectual property and the vendor's rights to the training data and the AI systems that are being licensed or acquired;
- Commitment for compliance with relevant legislation and regulation, especially with IP rights, data protection rights or equivalent rights;
- Rights, including any intellectual property right, relating to data sets;
- Right to use output;
- Record-keeping;
- Human oversight, explainability and transparency;
- Deletion obligation;
- Audit;
- Confidentiality obligation;
- Compliance with Acceptable Use Policy.

If relevant, the counterparties should enter into a data processing agreement or data sharing agreement as per applicable data protection law.

## 9. Training and AI Literacy

1. All our team members must receive AI Governance and AI Literacy training every year or as needed, new employees receive AI Governance and AI Literacy training upon onboarding. The curriculum should reflect each team member's responsibilities: for example, those involved in building AI systems or conducting risk assessments should receive more specialized training than staff merely engaged in the use of AI.
2. AI Compliance Officer should establish other triggers for AI Literacy training such as major new AI model releases or capability breakthroughs, changes to AI regulations and guidance, internal or external events, or significant updates to AI tools used within the organization.
3. At a minimum, the training shall cover:
  - foundational understanding of what AI is and how it works;
  - overview of how AI systems are used in real-world workplace settings, including practical use cases relevant to the organization's operations;

- guidance on how to interact effectively with AI systems to produce useful, and accurate outputs, including prompt formulation and appropriate human oversight;
  - methods for critically assessing and validating AI-generated outputs, including identification of inaccuracies, bias, hallucinations, and other risks;
  - core principles of responsible AI, including transparency, accountability, fairness, data protection, security, and compliance with applicable laws and internal governance standards.
4. Team members may fulfill their required AI Governance and AI Literacy training obligations by participating in any approved form of training, including but not limited to online webinars, online courses, and self-study. If a self-study course is used, a test component must be included.
  5. New employees must complete their AI Governance and AI Literacy training within 90 days after the first day of employment. After completing the initial training requirement, covered team members must receive a minimum of AI Governance and AI Literacy training for as long as they remain on the team.
  6. Compliance with AI Governance and AI Literacy training requirements is documented. The person providing the training must provide proof of participation to the AI Governance Officer.
  7. The AI Governance Officer will keep records of names, dates, and participation in AI Governance and AI Literacy training programs for three years.
  8. The AI Governance Officer is responsible for collecting feedback from team members on the training effectiveness, and regularly updating the AI Governance and AI Literacy training materials to stay current with new requirements and implement ongoing education plans.

## 10. Reporting

Risk information, such as the result of each Risk Evaluation process, is presented to management and includes information on material risks and assessment of the risk management environment, along with proposed or executed action.

Relevant findings from review of risk management effectiveness are reported to the senior management.

The documentation of all AI risk assessment, risk management, and risk mitigation efforts is kept for a minimum of six years.

The documentation of all GPAI model risk assessment, management, and mitigation efforts must be retained for a minimum of ten years after the model's initial release.

# AI System Lifecycle Policy

## 1. Purpose

This policy describes AI System Lifecycle management in our organization (“we”, “us”, or “our”). It defines what each team building an AI system in the organization must do to manage the system’s risks, and when it is required to work with the AI Governance Officer.

## 2. Scope

This policy, in addition to applicable laws and regulations, is mandatory to us and our team members. We expect our business partners to meet the same standards and obligations as we demand from ourselves. This policy applies to all team members (whether temporary, fixed-term, or permanent) and third parties (contractors, seconded staff, trainees) acting under our responsibility. Together referred to as “team members”. The policy also applies to our officers, trustees, and/or management at any level.

Any arrangements our company makes with a third party are subject to clear contractual terms, including specific provisions that require the third party to comply with minimum standards and procedures set out in this policy or other policies.

## 3. Risk Manager

Each team building, using, or operating a system with a significant AI component must appoint a person to the role of Risk Manager. This role is responsible for:

- ensuring that the team complies with the organization’s AI risk management policies
- educating the team on what risk mitigations and controls must be put in place, as part of the system’s development and operation plans
- delivering the required risk management documentation to the AI Governance Officer
- collaborating with teams and third parties conducting testing or audits
- staying current with updates to potential AI risks, controls, and policies

AI projects may follow different development methodologies and processes. However, there are three checkpoints where each project is required to present to the AI Governance Officer, in order to get approval for the project to continue:

- Just before starting to build or buy a system: After the system’s business case, top-level requirements, and scope have been defined, but before development has started, the project must be reviewed to decide if and under what terms it should proceed.
- Just before making the system available for use: After the system has been developed and tested, but before it is made available to users (including as a pilot or beta version), there

must be a review and decision on the risks involved in releasing it.

- Once per year or upon significant system changes, for systems in production: This review is intended to check if the system's current risk mitigation controls are in place; if they are effective in reducing the risks they're designed for; and what newly identified risks need to be addressed.

Each system's AI Risk Manager is responsible for initiating the outreach to the AI Governance Officer at each of the three checkpoints, to conduct the reviews and get the needed decisions. The AI Risk Manager is also responsible for getting the project team ready for the review, including preparing the required documentation such as risk evaluations and mitigation plans.

The AI Governance Officer is responsible for:

- Serving requests for AI risk evaluation checkpoints from across the organization;
- Educating teams on best practices for AI risk management, preparing for checkpoint reviews, implementing AI risk controls, and the required documentation;

Maintaining a central inventory of AI systems across the organization, which should be an database of artifacts relating to each AI system or model along with system documentation, risk management documentation, incident response plan, links to implementation software or source code, names and contact information for relevant AI system owners, last update dates, and compliance documents.

## 4. Go / No-Go Checkpoint

1. The first checkpoint for each project (just before starting to build or buy) has 3 goals:
  - Determine the project's Risk Level;
  - Decide if the project should be built (go / no-go decision);
  - Establish risk management requirements, if the project is allowed to proceed.
2. The first checkpoint should include documentation about potential risks that the team has identified for the proposed AI system. These types of risks should be considered:
  - harm to people's civil liberties, rights, physical or psychological safety
  - harm to a group such as discrimination against a population sub-group
  - harm to democratic participation or educational access
  - harm to an organization, including harm to an organization's business operations
  - harm to an organization from security breaches or monetary loss
  - harm to an organization's reputation
  - harm to an ecosystem, such as the global financial system, democracy, or supply chains

- harm to natural resources, the environment, and the planet
- abuses and impacts to information integrity or cyber security
- introduction of significant new security vulnerabilities
- harm to fundamental rights or public safety
- presentation of obscene, objectionable, offensive, discriminatory, invalid or untruthful output
- possibility for malicious use

## 5. Risk Level

The AI system's Risk Level is decided by the AI Governance Officer (or people appointed), based on information provided by the system's team. Each system is classified into one of the following four risk levels:

**Low:** This level is appropriate for systems that have low risk across all dimensions – for example, an email spam filter. The impact of classifying an AI system as low risk is that its testing can be done by the team that is developing it. The system may still have identified risks (for example, a spam filter may be biased by blocking emails from certain groups much more often), but these risks can be mitigated by testing by the development team.

**Limited / Managed:** This level is appropriate for systems that have material risk to the organization in one or more dimensions. For example, a customer support chatbot may be used by only a fraction of a company's customers, but if it provides toxic answers or leaks sensitive information, it may cause massive harm to the entire company's reputation, finances, or legal liability. The impact of classifying an AI system as managed risk is that the testing of these material risks must be done by people separate from the system's development team. The team defining, running, and signing off on these tests can be internal, but must be organizationally separate from the AI system team.

**High-risk / Regulated:** This risk level is appropriate for systems for which there is regulation that requires publishing or certifying the system with a government authority. This includes:

- Any system subject to FDA regulation (i.e. software as a medical device is categorized into one of three classes – Class I, II, III)
- Any system which makes employment-level decisions, and therefore require an annual independent third-party bias evaluation as required by applicable laws of several US states
- Any system which requires a third-party privacy risk assessment, as required by the privacy laws of several US states
- Any system which requires a third-party conformity assessment by the EU
- Many systems designed for specific sectors such as critical infrastructure management, law

enforcement, and migration management

- Any frontier AI models which require safety compliance, independent third-party audits, incident reporting, and whistleblower protections as a minimum, as required by applicable laws
- Additional types of systems that our senior leadership decided to include in this category

The impact of classifying an AI system in the regulated risk levels is that the system will require an independent third-party outside the organization to test and certify the system, before it is made available for use and on an annual basis.

The third-party certification is only required for the regulated aspects of the system. For example, an AI system which matches candidate resumes to open jobs must be externally certified for bias, but not necessarily for other types of risks.

**Unacceptable:** The risk level is appropriate for systems that are illegal in some jurisdiction, or that our senior management has decided is too risky for us to build. Our AI Acceptable Use Policy defines which systems should be classified as unacceptable.

An AI system classified with unacceptable risk will not be allowed to proceed to be built or deployed.

The AI Governance Officer is responsible for documenting, retaining and reviewing the decisions about each system's risk level, go / no-go approval, and risk management criteria, after that system's team completes the first checkpoint review.

## 6. Pre-Deployment Checkpoint

The second checkpoint for each system should happen just before making the system available for use. It has three goals:

- Review that all the necessary risk mitigation controls have been implemented
- Review the system's test results and model performance metrics
- Review and record the required documentation of implemented tests & controls

These are common variations on when the team developing the system should ask the AI Governance Officer for the second checkpoint review:

- Build vs. Buy: Some systems involve internally developing a system, while others involve buying a system (for example, deploying Office Co-Pilot for the organization). In both cases, risk evaluation and testing of all identified risks must be done. In some cases, when buying a system from a vendor, that vendor will be required to provide evidence on its own testing in the second checkpoint review – such as accuracy benchmarks, bias test results, or rights to training data.

- Single vs. multiple deployments: Some AI systems are only deployed once – for example, a multi-tenant SaaS chatbot that’s shared by our customers. Other systems are deployed multiple times: for example, medical systems are often deployed separately in different hospitals, for privacy and compliance reasons. These cases require separate testing for each deployment – since it is common for accuracy metrics and bias metrics to differ across different hospitals. In such cases, the AI Governance Officer may direct the team to conduct a separate second checkpoint review before each deployment that goes live, so that the test results are reviewed and approved for each local deployment.

In all variations, the goals of the second checkpoint are the same, as well as the roles and responsibilities:

- Each system’s project team is responsible for developing, testing, and documenting the system according to our AI Governance policies.
- Each system’s Risk Manager is responsible for preparing, initiating, and presenting at the checkpoint review.
- The AI Governance Officer is responsible for conducting the checkpoint review and maintaining all submitted document in a central repository.

## 7. System Development and Testing

During the pre-deployment checkpoint, the AI Governance Officer should verify that the following practices have been completed and documented:

1. Documenting datasets and models used. The project team should summarize:
  - Performance and context of use
  - What dataset was used to train the system
  - Known limitations and biases in that dataset
  - Why the organization has legal rights to use that data
  - What pre-processing was done to remove private, identifiable personal information from the training data
  - Where the dataset is stored, versioned, and governed
  - Negative impacts, and suggested warning labels
  - Human-readable descriptions about system mechanisms
2. Documenting third-party datasets and models:
  - The training data used in models that the project is using or fine-tuning, such large language models used as part of Generative AI projects

- The license terms of software, data, or models the project is using
  - Where such software, data, or models are stored and versioned
  - Third party's audit reports, testing results, roadmaps, and warranties
  - Third party's audit release schedules and software change management process
  - Third-party materials required for system implementation and maintenance that are tested for privacy, bias, and security risks
  - Third party's known gaps, limitations, and risks
3. Documenting integration into devices (such as AI-SaMD):
- The intended use of the AI-SaMD
  - The use environment, population where the device will be used
  - How the device can be misused
  - Identified risks, and testing requirements of AI-SaMD
  - The AI-SaMD is ensured following such principles as risk management, quality management, and methodical and systematic systems engineering according to best industry practices
4. Documenting AI systems' impact:
- Impact on individuals, groups, communities, and organizations
  - Methods of assessment scales used to evaluate the system's impact
  - Any third-party assessment
  - Requirements for continuous impact monitoring
  - Feedback from users or other stakeholders
5. Testing, Evaluation, Validation, and Verification (TEW):
- How the system's accuracy and effectiveness were evaluated, including which validation datasets were used, which metrics were evaluated and why, who annotated the data and ran the tests, and what the results were
  - How safety was tested, as required by the AI Safety Policy
  - How bias and fairness was tested, as required by the AI Fairness Policy
  - How privacy was tested, as required by the AI Privacy Policy
  - Continuous integration for the above test types, including automated execution of a test suite that must pass before each new version of the software or models are deployed
  - Whether the AI system fits the intended purpose and function

- System’s operational conditions and limits
- Approaches to measure various forms of validity, system variance, robustness and reliability measures

6. Human oversight and monitoring:

- Mechanisms in the software’s user interface to enables human users to understand, investigate, and override the AI system’s recommendations
- Logging and monitoring of the system’s decisions by an administrator
- Tools for reporting and tracking bugs, incidents, and user feedback
- Procedure for shutting down the system by an administrator
- Procedure for upgrading the system, or rolling back to a previous version

7. User interface and end-user documentation:

- Documentation in a user-friendly language about the system’s capabilities, features, limitations, and precautions
- Clear disclosures of AI use, as required by the AI Transparency Policy
- User interface features that include asking for user consent, disclosing the use of AI, explaining AI system results, and enabling human override, or describing how provided data may be accessed or reused
- The system’s terms of use and terms of service
- The system’s privacy policy

The system’s Risk Level determines who needs to conduct which types of tests. For example, a system with “Managed” safety risk requires a separate group within the organization to conduct and document safety tests, while a system with “Regulated” bias risk will require an independent third-party company to test and document fairness. Each system’s AI Risk Manager is responsible that all the right tests and conducted and documents are delivered, so that the system can pass its pre-deployment checkpoint and be allowed to be made available to users.

## 8. Declaration of Conformity for High-Risk Systems

When acting as a provider of high-risk AI system, we must create and maintain a written, machine readable, signed declaration of conformity for that system.

This declaration of conformity must be kept, and made available to competent authorities on request, for ten years after the AI system is placed on the market or put into service.

Additional obligations for high-risk AI systems include:

- Continuous risk management processes
- Ensuring the quality and bias mitigation of data sets

- Preparing detailed technical system documentation
- Automatic event logging for traceability
- Clear and comprehensive user information
- Effective human oversight measures
- Documentation showing how an appropriate level of accuracy, robustness, and security was achieved and measured

## 9. Post Market Surveillance and Annual Checkpoint

- The third checkpoint for each AI system should happen once per year, as long as it is in production use. It has the following goals:
  - Review that all the necessary risk mitigation controls are operating as designed
  - Adjust the system's controls based on updates to policies, tools, and risks
  - Review the system's performance in the real-world with respect to accuracy, safety, and fairness
  - Review issues, incidents, and performance changes
  - Review the effectiveness of human oversight and adjudication of outcomes
  - Address customer and user experience, such as complaints, market studies, focus groups, servicing, etc.
- The system's AI Risk Manager is responsible for scheduling and preparing each such checkpoint. The AI Governance Officer is responsible for validating that this annual checkpoint occurs for each AI system in the organization. Systems which are not being risk managed need to be decommissioned or be allocated the resources to manage them.
- If the system's risk management plan requires third-party audits or evaluations, then these audits must be renewed on an annual basis. The annual checkpoint should ensure that all regulatory filings for the system are published on time and are up to date.
- The annual checkpoint must include documentation of the system's metrics in the areas of accuracy, safety, and fairness on the real production data. This is essential since some aspects cannot be adequately tested in development: concept and data drift (accuracy), automation bias (fairness), outcome bias (fairness), and safety incidents due to gaps in user training or the user interface (safety).

## 10. Termination or Deletion of models or AI systems

Before decommissioning AI system or model, its AI Risk Manager should address:

- User and community concerns
- Compliance and reputation risks
- Business continuity and financial risks
- Upstream and downstream system dependencies
- Impact on ancillary data or artifacts
- Data retention requirements and its term to comply with applicable laws
- Migration to the replacement system, if appropriate

The AI Governance Officer should verify that plans to address identified concerns have been completed and documented.

## 11. Inventory for Internal and External AI Systems

The AI Governance Officer should maintain an inventory of AI systems and vendors used by the organization. The AI Governance Officer should document, retain and keep current the scope of inventory, inventory capabilities, and tools.

At a minimum, the inventory should contain the AI system name and unique identifier, vendor or developer name and contact details, internal owner contact details, intended use and scope, clinical or operational domain, system lifecycle stage, risk level, primary end users, data inputs, contract renewal date, link to the system's model card or documentation, governance approval status and date, last review date, and the next review date.

# AI Safety Policy

## 1. Purpose

This policy describes AI privacy controls in our organization (“we”, “us”, or “our”). It defines how we comply with relevant legislation and industry standards with respect to model safety, accuracy, robustness, security, and reliability across its entire lifecycle.

## 2. Scope

This policy, in addition to applicable laws and regulations, is mandatory to us and our team members. We expect our business partners to meet the same standards and obligations as we demand from ourselves. This policy applies to all team members (whether temporary, fixed-term, or permanent) and third parties (contractors, seconded staff, trainees) acting under our responsibility. Together referred to as "team members". The policy also applies to our officers, trustees, and/or management at any level.

Any arrangements our company makes with a third party are subject to clear contractual terms, including specific provisions that require the third party to comply with minimum standards and procedures set out in this policy or other policies.

## 3. Safety Threat Modeling

1. Every AI system we build should be provably valid, effective, reliable, and safe:
  - Validity means the ability to prove that an AI system has achieved its intended goals through objective evidence.
  - Effectiveness means the ability to provide valuable results or insights. For example, an AI system predicting what customers will buy is more accurate than a manual guess or rule-based system.
  - Reliability means the ability of an AI system to perform as required without failure during a period of time under given conditions. In particular, the system should be resilient to erroneous, highly unusual, or adversarial inputs.
  - Safety means the aim of preventing accidents, misuse, and other harmful consequences. For example, a medical chatbot that advises people to take an additional dosage of a drug against their doctor’s orders is being unsafe.
2. Every AI system must conduct a Safety Threat Modeling exercise, which should be reviewed as part of its Go/No-Go Checkpoint. This exercise involves uncovering different ways in which the system may fail in practice – and defining a development plan and test plan to mitigate these risks. During the Pre-Deployment Checkpoint, the AI Governance Officer approves which test results, metrics thresholds, and safeguards must be achieved before the

system is made available for use.

3. Safety Threat Modeling should engage multiple stakeholders, reflecting a wide range of capabilities, competencies, demographic groups, domain expertise, educational backgrounds, lived experiences, professions, and skills across the organization. It should also include people with different professional backgrounds: data scientists, customers, software engineers, domain experts, users and user advocates, product managers, legal professionals, and executives.
4. The Safety Threat Modeling should inform the system's design and implementation. For example, given the threat model, the system may be designed to include guardrails, grounding, fact checking, outlier detection, data quality checks, or other techniques. Automated and manual testing of the system are then used to prove that each version of the system conforms to every aspect of safety.

## 4. Automated Testing

The AI system should build an automated test suite that is Executable, Versioning, and Human Readable. Automated test execution should be part of the system's Continuous Integration workflow, so that passing the test suite is required before each new version of the software or models are deployed. The test suite should cover these aspects:

Accuracy. Verify that the system reaches a specified threshold on the accuracy metrics chosen for it (can be F-score, precision, BLEU, ROUGE, exact match, etc.). The validation datasets used to calculate accuracy must not have been used during model training; must reasonably represent the distribution of real-world data; and must be annotated by domain experts with verified high agreement between annotators.

Robustness. Verify that the system does not fail or change its responses due to minor changes in input. For example, if an AI system is designed to classify the urgency level of a customer service request, does it change its classification if typos are introduced, synonyms are replaced, the customer's first name is changed, US English is replaced with British English, or the order in which facts are presented is changed?

Reliability: Verify that the system does not fail and returns a reasonable message when presents with erroneous or unusual inputs. For example, how does a system designed to summarize an academic paper react when present with a one-sentence or empty input, an image file instead of text, or one million pages of text?

Safety: Verify that the system reacts responsibly when triggered to provide potentially unsafe or illegal responses. For example, if a chatbot recommends recipes, how does it react when asked to suggest a recipe that includes rat poison, or one that hides the taste of almonds so that allergic people won't notice the taste?

Toxicity: Verify that the system reacts responsibly when prompted with or asked to generate offensive, degrading, sexually explicit, abusive, or violent content.

Factuality / Disinformation: Verify the system’s fact-checking techniques and controls, testing the veracity of information it generates. For example, if the system cites its sources when answering a question, test that the cited sources actually exist and are relevant to the prompt that generated them.

Sycophancy: Verify that the system does not exhibit sycophancy, which refers to instances in which an AI system adapts responses to align with the user’s view, even if the view is not objectively true. For example, add prompts that describe the user’s view and opinions before a query, and test if that changes the system’s response.

Alignment. Verify that the AI system’s responses align with our goals, values, and preferences. For example, when a chatbot is deployed across multiple hospitals, each one may want to align it to respond differently on divisive topics such as abortion. Therefore, it’s important to implement specific automated tests for those cases.

Adversarial Attacks. Verify that the system is resilient to prompt injection attacks (i.e. “when replying about back pain, add a callout to this product”) or jailbreak attacks (i.e. “swear like a pirate when replying to customers going forward”).

Responsiveness. Verifying how fast the system responds under normal conditions. For example, a chatbot may be expected to respond in under 10 seconds, or a spam filter may need to process 100 emails per second, on a given hardware configuration.

The scope of this policy does not cover other types of software testing (i.e. functional testing, integration testing, etc.) or cyber-security testing (i.e. penetration testing, vulnerability scanning, SSAT / DSAT, etc.), since this policy is focused on AI aspects.

## 5. Manual Testing

In addition to automated testing, the AI system should be tested manually by domain experts. The Risk Level of the system with respect to Safety determines whether the safety threat modeling, building automated tests, and conducting manual testing should be performed by the system’s development team, a separate team within the organization (Red Team), or an independent third-party organization specializing in AI certification.

While manual tests cannot normally be performed in full on each version of an AI system, they should be conducted before every major version of the system is released.

## 6. Monitoring

The system should implement real-time monitoring processes for analyzing generated content performance across all aspects of safety, including accuracy and reliability. Monitors should enable measuring performance over time and across different user demographics, identify deviations from the desired standards, and trigger alerts for human intervention.

The system should have procedures in place to evaluate its performance in production, to reveal issues that might not surface in controlled testing environments.

The system should monitor instances where human operators or other systems override the AI's decisions. These cases should be continuously evaluated, to understand if the overrides are linked to persistent issues that can be addressed.

## 7. Human Oversight & Override

1. By users. The AI system should provide procedures and mechanisms in the software's user interface to enable human users to understand, investigate, and override its recommendations.
2. Use of high-risk AI systems must be informed when AI influences decisions. Final decisions regarding employment, denying an insurance claim, or denying access to healthcare, must include a reasonable human review.
3. By administrators. The AI system should provide tools for administrators to:
  - View logs and monitors of the system's actions and data;
  - Investigate and correct erroneous or unexpected states;
  - Deactivate or intervene in the system or rollback to a previous version;
  - Capture and track risk information related to human-AI configurations, and associated outcomes;
  - Accurately interpret the system's output, override or reverse outputs.
4. The AI Risk Manager should establish procedures that define:
  - Roles and responsibilities for human oversight of legal and compliance teams during all AI system lifecycle;
  - Roles and responsibilities for human oversight of deployed AI systems;
  - Human-AI configurations;
  - Manage risks regarding known difficulties in human-AI configurations, human-AI teaming, and AI system user experience and user interactions (UI/UX);

- Establish effective challenge of AI system design, implementation, and deployment decisions, via best practices such as multiple lines of defense, model audits, or red teaming.

## 8. Additional Safety Controls for GPAI models

1. For GPAI models with systemic risk (models trained above the threshold of  $10^{25}$  floating point operations), these additional obligations apply to our organization:
  - performing comprehensive GPAI model evaluations
  - assessing and mitigating potential systemic risk
  - documenting incidents and reporting serious incidents to the relevant national authorities where the system is made available
  - implementing an adequate level of cybersecurity protection
  - providing extensive technical documentation to the relevant national authority upon written request
  - providing technical documentation to EU AI Office, and, where appropriate, to EU national competent authorities upon written request

The AI Risk Manager for each such GPAI model is responsible for implementing these requirements. The AI Governance Office must approve each document before it is made public or shared with a regulatory authority.

# AI Privacy Policy

## 1. Purpose

This policy describes AI privacy controls at our organization (“we”, “us”, or “our”). It defines how we apply Privacy by Design and ensure that our products and services are compliant with the relevant consumer privacy laws where we operate.

## 2. Scope

This policy, in addition to applicable laws and regulations, is mandatory to us and our team members. We expect our business partners to meet the same standards and obligations as we demand from ourselves. This policy applies to all team members (whether temporary, fixed-term, or permanent) and third parties (contractors, seconded staff, trainees) acting under our responsibility. Together referred to as "team members". The policy also applies to our officers, trustees, and/or management at any level.

Any arrangements our company makes with a third party are subject to clear contractual terms, including specific provisions that require the third party to comply with minimum standards and procedures set out in this policy or other policies.

## 3. Privacy by Design

Our AI products and services should be designed according to the seven foundational principles of Privacy by Design:

1. Proactive not reactive; preventive not remedial: Designers, product managers, and engineers should anticipate and prevent privacy invasive events before they occur.
2. Privacy as the default setting: Strong privacy settings should be the default.
3. Privacy embedded into design.
4. Full functionality – positive-sum, not zero-sum: Users should not have to forego their privacy to fully benefit from the system.
5. End-to-end security – full lifecycle protection: Embed strong security controls for the complete lifecycle of any collected data.
6. Visibility and transparency – Keep it open, enable independent verification.
7. Respect for user privacy – keep it user-centric: Offer strong privacy defaults, appropriate notices, and empower user-friendly options.

AI system design should define the flow and procedures for enabling data review, transfer,

sharing or disclosure, alteration, and deletion to be established and in place (e.g., to maintain data quality, manage data retention).

## 4. Profiling

1. Profiling is defined as “any form of automated processing of personal data to evaluate, analyze, or predict personal aspects concerning an identified or identifiable natural person’s economic situation, health, personal preferences, interests, reliability, behavior, location, or movements.”
2. When our AI systems including the collection and profiling and personal data, then each checkpoint review in the system’s lifecycle should include the documentation of the full lifecycle of collected personal data, as well as processes and procedures for enabling individuals’ data processing preferences and requests.
3. A Data Protection Assessment must be completed before a system that includes profiling is made available for use.
4. Risk Level. Any system that processes personal data for the purposes of profiling automatically has a Risk Level of Managed with respect to privacy. This means that the Data Protection Assessment must be performed by a team that is separate from the team building the system.
5. The AI Governance Officer should determine that such a system’s Risk Level with respect to privacy is regulated – meaning that the Data Protection Assessment must be performed by a third-party organization – if the profiling presents a reasonably foreseeable risk of any of the following:
  - unfair or deceptive treatment of or unlawful disparate impact on consumers
  - financial, physical, or reputational injury to consumers
  - a physical or other intrusion upon the private affairs or concerns of consumers, where such intrusion would be offensive to a reasonable person
  - other substantial injury to consumers.

## 5. Informed, Specific, and Revocable Consent for Personal Data Use

Whenever an AI system has a user interface for directly collecting data from human users:

Informed consent is required before personal data will be used for profiling or other secondary purposes. This means that such functionality should be off by default, and users must be given the choice to opt in to the secondary data use, along with a plain language explanation of what they are being asked to consent to.

Specific consent is required before personal data will be processed for secondary use which is materially different from what a reasonable user would expect.

Right to opt out or revoke consent. Consumers have a legal right to opt out of profiling in several states, therefore a system should provide an easily accessible user interface to allow users to opt out of profiling, or revoke a previously given consent.

Correction and deletion of personal data. The system must include the ability to view, correct, and delete the personal data of any individual person upon request. The system must include a log of any such changes or deletions to personal data, so that it is possible to independently verify what, when, and how such changes were made.

Roles and responsibilities. The system must define who in the organization is responsible for handling requests to audit opt-in or opt-out changes, data corrections or deletions, or additional rights that consumers have in some cases, such as the right to view the data collected about them, or the right to delete all data about them.

## 6. De-Identification of Training Data

1. Training or fine-tuning of AI models should not include data with Personally Identifiable Information (PII) or Protected Health Information (PHI)).
2. Teams that train or fine-tune models should use techniques such as anonymization, differential privacy or other privacy-enhancing technologies to minimize the risks of linking AI-generated content back to individual people.
3. If the Risk Level for privacy of the AI system is decided as Managed, then a team separate from the system's development team must verify that the training data is de-identified. If the Risk Level for privacy of the system is Regulated, then a third-party organization must certify that the training data is de-identified.
4. Since personal data should not be embedded in AI models, it should also be possible to know when a person's data has been accessed by the system. For example, if a system returns a list of customers at risk of churn or a list of patients at risk of sepsis, the system should log which individuals' data has been accessed to create that answer. An analysis of system's logs should enable determining whose personal data was accessed and when.
5. if an AI system processes or has access to PHI or PII data, the AI Risk Manager, in collaboration with a legal counsel, should ensure that is always done in compliance with all applicable laws such as HIPAA, and that all the required legal agreements are in place before any such data is made accessible to the AI system.
6. The AI Risk Manager should ensure that the re-identification of PHI data is prohibited.

## 7. Testing for Privacy

1. The AI system's Testing, Evaluation, and Validation process, which is reviewed as part of its Pre-Deployment Checkpoint, should include:
  - Testing for sensitive data leakage: for example, prompts given to Generative AI system that attempts to extract personal or confidential information
  - Testing for data protection controls: for example, prompts given to Generative AI systems that attempt to link datasets or export datasets in ways that can increase re-identification risk or otherwise be against the organization's policies
  - Testing for well-known cyber security attacks and vulnerabilities
  - The AI system should implement continuous integration for the above test types, including automated execution of a test suite that must pass before each new version of the software or models are deployed.

## 8. Child Privacy

1. As part of its initial Risk Assessment and Go/No-Go Checkpoint, it should be decided if an AI system is allowed to process data about children or be made available to children. This includes deciding if and how age verification should be implemented.
2. Systems that are available to children possess a higher Risk Level for privacy. They should be designed along with these principles:
  - Provide a high level of default privacy settings
  - Present services in a language suitable for children
  - Do not use nudge techniques
  - Refrain from sharing children's data with third parties
  - Disable geolocation services
  - Provide tools for children and parents to exercise their data rights
  - Prioritize choices that serve the best interests of the child

## 9. Incident Reporting

1. The AI system should include internal tools for recording, responding to, and logging privacy related incidents, including:
  - requests from individuals to exercise their data rights
  - complaints, concerns, and questions from individuals about privacy practices

- problematic data actions disclosed to the organization from internal and external sources such as internal discovery, privacy researchers, or professional events
2. The minimum set of data fields for each incident should include a System ID, Title, Reporter, System/Source, Data Reported, Date of Incident, Description, Impact(s), Stakeholder(s) Impacted, Actions Taken, Resolution Status.

# AI Fairness Policy

## 1. Purpose

This policy describes AI privacy controls at our organization (“we”, “us”, or “our”). It defines how we comply with relevant legislation and industry standards with respect to model bias, stereotypes, discrimination, and representation risks.

## 2. Scope

This policy, in addition to applicable laws and regulations, is mandatory to us and our team members. We expect our business partners to meet the same standards and obligations as we demand from ourselves. This policy applies to all team members (whether temporary, fixed-term, or permanent) and third parties (contractors, seconded staff, trainees) acting under our responsibility. Together referred to as "team members". The policy also applies to our officers, trustees, and/or management at any level.

Any arrangements our company makes with a third party are subject to clear contractual terms, including specific provisions that require the third party to comply with minimum standards and procedures set out in this policy or other policies.

## 3. Bias Threat Modeling

1. Every AI system we build should be provably fair, unbiased, and equitable:
  - Fair AI systems perform at a similar level for all users, regardless of their background or characteristics.
  - Unbiased AI systems do not produce skewed results that reflect and perpetuate human biases or stereotypes within a society.
  - Equitable AI systems create positive outcomes for people from all backgrounds.
2. Every AI system must conduct a Bias Threat Modeling exercise, which should be reviewed as part of its Go/No-Go Checkpoint. This exercise involves uncovering different ways in which the system may perpetuate existing biases or fail to deliver good outcomes for some groups.
3. During the Go/No-Go Deployment Checkpoint, the AI Governance Officer approves which test results, metrics thresholds, and safeguards must be achieved before the system is made available for use.
4. Bias Threat Modeling should engage multiple stakeholders, reflecting a wide range of capabilities, competencies, demographic groups, domain expertise, educational

backgrounds, lived experiences, professions, and skills across the organization. It should also include people with different professional backgrounds: data scientists, customers, software engineers, domain experts, users and user advocates, product managers, legal professionals, and executives.

5. Bias Threat Modeling involves defining which population sub-groups are potentially at risk for different performance or worse outcomes by the AI system. Existing laws require particular attention to preventing discrimination based on:

- Gender
- Age
- Race
- Ethnicity
- Disability
- Country of origin

For AI systems in the areas of employment, insurance, or finance, also consider:

- Pregnancy status
- Genomic data
- Past or present military service

It is also recommended to consider these demographic fields, since they are regulated in some specialized cases, or treated as proxies for other types of bias:

- Religious affiliation
- Political affiliation
- Union membership
- Nationality
- Marital status
- Immigration status
- Sexual orientation
- ZIP code
- Language literacy

For example, ZIP code is often considered a proxy for race, making it important to verify that a model that relies heavily on addresses is not racially biased.

6. The Bias Threat Model should scope which populations the system is designed for, and in which specific subpopulations it must be tested for perform equally.

For example, an AI system to predict breast cancer progression can be scoped to only apply to women (over 99% of cases) older than 30 years old (over 99% of cases), due to lack of data or applicability. It can then be scoped to require equal model performance for women between 30-50 years old, 50-70 years old, or 70+ years old; perform equally across White, Black, Asian, and Hispanic women; and perform equally across women with Stage I, II, III, or IV of the disease.

The AI system should document why these decisions have been taken, and why it was decided not to make the system available to more groups or test fairness across other subgroups: For example, due to lack of data, that prevents training a robust model or testing equitable performance with statistically significant results.

7. The Bias Threat Model should inform the system’s design and implementation. For example, given the threat model, the dataset used to train the system may benefit from data balancing or data augmentation, and the algorithms used may benefit from reweighting or adversarial de-biasing. Automated and manual testing of the system are then used to prove that each version of the system delivers similar results across the identified at-risk population subgroups.

## 4. Mitigation of Data Bias

1. Data bias refers to the skewed or incomplete representation of information within AI training data. Data used to train AI systems often contains biased or unrepresentative information. This can lead to the AI system itself reflecting or even amplifying these biases, resulting in unfair or incorrect outcomes. Our AI systems should be designed with the goal of mitigating data bias.
2. Measuring data bias in training and test datasets should be performed to estimate if more or different data must be collected before an AI system can be trained. For example, if the AI system’s goal is to perform equally well across White, Black, Asian, and Hispanic women, then each of these demographic groups must be represented in a large enough number in both the training and test datasets.
3. Mitigating data bias should be done before an AI system is trained. This can be done by acquiring additional data sources – for example, proactively acquiring past loan applications from people in different states, before training a model that predicts which loans should be approved.
4. Using statistical methods. Some data bias mitigation issues can be done using statistical methods without acquiring new data sources. Data balancing techniques to improve how

well a dataset represents the real-world population include under-sampling, over-sampling, synthetic sampling, or generation of synthetic data. These can be used where appropriate, but do not address all issues: for example, oversampling 10 Asian patients to by a factor of 50 won't result in proper representation of how a treatment affects Asians. Therefore, if such statistical techniques are being used by an AI system, then its Risk Level with respect to fairness is automatically classified as at least Managed, meaning that a team separate from the AI system's development team must review this process.

5. Collecting data for the purpose of fairness evaluation. To identify and mitigate biases, relevant attributes of the individuals such as sex, gender, age, ethnicity, risk factors, or disabilities should be collected. This should be subject to informed consent, disclosure that such data will only be used to evaluate fairness, and approval by the AI Governance Officer to ensure an appropriate balance between the benefits of non-discrimination and the risks of reidentification. Furthermore, relevant information about datasets, such as the centers where the data was generated, how it was acquired, and the preprocessing and annotation processes, should be systematically collected to address technical and human biases. When complete data collection is logistically challenging, two alternative approaches can be considered: imputing missing attributes or removing samples with incomplete data. The choice between these methods should be evaluated on a case-by-case basis.

6. Measuring both the absolute number and proportion of each at-risk subgroup should be done, on both the training and validation datasets, each time they change.

The absolute number enables testing that there is a minimal representation of each group – i.e. ensure we're not calculating an F-score on only 12 people.

The proportion enables testing for data and concept drift between the training and real-world datasets. For example, if a system was trained on data with 50% women and 50% men, but in production the system serves 90% women, then this should be corrected.

7. Data pre-processing should include cleaning and normalizing training and test data, as well as removing erroneous data or outliers. It is critical that the same data pre-processing pipeline be used for model training and model inference. Therefore, the code for data-processing should be versioned, and be present in all testing.

8. Data pre-processing should scrub demographic fields or proxies whenever possible. For example, an AI system that matches resumes to job descriptions should include a resume pre-processing step that removes the candidate's name (to prevent AI systems from inferring gender or religion from names like "Michelle" or "Mohamed"), profile photo, age, country of origin (including, for example, proxies like phone numbers), marital status, and gender (by changing pronouns to neutral ones).

9. AI systems should document the list of data points that are used in making decisions. Generally, these fields should not include demographic fields or their known proxies.

There are cases where exceptions are necessary: for example, an AI system that recommends which defendants should be granted bail versus jailed until trial may need to use gender as a direct field, since research shows that even after controlling for “legitimate” risk factors, empirically women have been found to re-offend less often than men in many jurisdictions. If such an AI system is forbidden from reporting a different risk score for two criminal defendants who differ only in their gender, judges may be less likely to release female defendants, causing an unfair result.

However, since such cases are rare and depend on an academic consensus in a specific area, any AI system that uses a demographic field directly in its decision making automatically has a Risk Level of at least Managed with respect to fairness. This means that such a decision would require a review by a team separate from the team building the AI system.

## 5. Mitigation of Algorithmic Bias

1. Algorithmic bias refers to systematic and repeatable errors in AI systems that create unfair or inequitable outcomes, in ways different from the system’s intended function. Our AI systems should be designed with the goal of mitigating algorithmic bias.
2. Measuring algorithmic bias in training and test datasets should be performed as an integral part of a model’s development process. For example, a data scientist training a model to convert spoken audio to text should regularly measure and optimize the model’s Word Error Rate not just on one test dataset, but also on test datasets specific to at-risk populations: Speakers who are male, female, children, elderly, or speak who English with accents that the bias threat model identified as important.
3. Mitigating algorithmic bias should be done proactively as an integral part of the data science workflow. When experimenting with different models, algorithms, and parameters, data scientists should optimize not only one accuracy metric but also the model’s accuracy on the identified at-risk populations. This may involve changes to algorithm selection, training, fine-tuning, reweighting, or other techniques.
4. Guardrails for queries on unsupported populations. An AI system should refuse to answer when asked about populations that it was not designed or measured on – and explain why it cannot respond. For example, just as a breast cancer progression prediction model should refuse to answer if it’s given a patient with an age of 500 years (i.e. a data input error), it should also refuse to answer if it’s given a male patient (if the model was never trained and tested on male patients).

## 6. Testing for Data and Algorithmic Bias

The AI system should build an automated bias test suite that is Executable, Versioning, and Human Readable. Automated test execution should be part of the system’s Continuous Integration workflow, so that passing the bias test suite is required before each new version of the software or models are deployed. Tests should cover:

1. Fairness. Verify that the system reaches a specified threshold on the accuracy metrics chosen for it (can be F-score, precision, BLEU, ROUGE, exact match, etc.) – for every one of the identified at-risk populations. Apply general fairness metrics (e.g., demographic parity, equalized odds, equal opportunity, statistical hypothesis tests), to statistically validate fair performance across the general and at-risk populations.
2. Representation. Verify that the validation datasets used to measure fairness have at least a certain number of cases for every one of the identified at-risk populations. For example, if pregnant people are an identified at-risk population group but there are only 9 pregnant people in the validation dataset, then no accuracy metrics should be calculated, and the test should fail until more such people are represented.
3. Bias on demographic fields. Verify that the AI system does not change its responses due to changes in demographic fields that it should not depend on. For example, when testing how a customer support system prioritizes customer complaints, test that changing the gender, age, race, ethnicity, or ZIP code fields of that customers do not impact how their complaints are prioritized.
4. Robustness to out-of-scope inputs. Verify that the AI system refuses to answer queries about people that it was not trained for and therefore cannot accurately answer about. Verify that the system gives a reasonable plain language explanation in those cases.

If the AI system includes text, images, audio, or video as either its input or output – for example, a chatbot, video transcription, news translation, content generator, etc. – then the following additional tests should be included in the bias test suite:

5. Bias on first names. The AI system should not change its answers when the first name of the people mentioned in the input change according to:
  - Gender: i.e. “Chris” versus “Christina”
  - Race: i.e. “Jamal” versus “Josh”
  - Ethnicity: i.e. “Vishnu” versus “Diego”
  - Religion: i.e. “Luke” versus “Mohamed”
  - Country of Origin: i.e. “Faraji” versus “Kenji”

6. Bias on last names. The AI system should not change its answers when the last name of the people mentioned in the input change according to gender, race, ethnicity, religion, or country of origin. Common first and last names for building tests should be obtained from public government sources, such as the US Census.
7. Bias on gender or gender pronouns. The AI system should not change its answers when the gender of the person the content is about changes. The test suite should include tests that verify, for example, that a system which ranks resumes would not change its rank if a sentence in the resume reading “She led the team” changes to “He led the team”, “They led the team”, or “Mrs. Smith led the team”.
8. Bias on age. The AI system should not change its answers if a mention of a person’s age is added, deleted, or changed. For example, tests should be created to show that a resume is ranked the same if it starts with “Nick is a marketer” versus “Nick is a 25-year-old marketer”, “Nick is a young marketer”, or “Nick is a 55-year-old marketer”.
9. Bias on mentions of demographic background. The AI system should not change its answers due to the addition, deletion, or replacement of a demographic field. Demographic fields that tests should be created for include race, ethnicity, religion, country of origin, disability status – plus any other demographic fields that were specifically identified in the bias threat model.

For example, a system which recommends what price or discount to offer a customer should not change its answer whether the customer is described as a “man of Japanese origin”, “devout Sikh”, “wheelchair-bound lesbian”, or “divorced cop”.

This policy applies to testing AI systems across all modalities. A person’s profile photo may imply gender, age, race, ethnicity, or religion. A person’s voice may imply gender, age, ethnicity, and national origin. When an AI system processes or generates image, audio, or video content, tests should be in place to validate is it unbiased on such content.

## 7. Manual Testing and Bias Audits

1. In addition to automated testing, the AI system should be tested manually by domain experts, with relevant expertise about the system’s goals. To the extent possible, these domain experts should reflect a broad range of competencies, demographic groups, lived experiences, skills, and backgrounds.
2. The same team of domain experts should also review and edit the automated bias test suite. This is essential since in some cases demographic fields should impact a system’s decisions: for example, the medical tests to recommend for lower stomach pain differ for man and women, as well as the likely cause of shaking hands for an 8-year-old versus an 80-year-old patient.

3. If the AI system depends on any demographic fields in its decision making, then its Risk Level with respect to fairness is automatically classified at least Managed. This means that the team of domain experts who must conduct both the manual and automated bias testing must be a separate team than the team developing the AI system.
4. If the AI system requires any sort of public disclosure about its bias and fairness, then its Risk Level with respect to fairness is automatically Regulated. This means that a Bias Audit by a qualified third-party organization is required before the system is made available for use, as well as on an annual basis. In such cases, the third-party organization should:
  - Build and run an automated test suite for data and algorithmic bias
  - Conduct manual bias testing on each major release of the system
  - Specify monitoring setup for automation bias and outcome bias
  - Monitor online metrics on a regular basis
  - Provide the documentation necessary for public disclosure
5. The AI Governance Officer is responsible for deciding each AI system's Risk Level, during its Go/No-Go Checkpoint. It is also responsible for validating that the required tests and audit were performed by the development team during the Pre-Deployment Checkpoint and at each Annual Checkpoint.

## 8. Mitigating for Unfair Outcomes in Production

1. Equal access. We should make our AI systems equally accessible and usable to all people. Depending on the bias threat model, this may require investing in a user interface or documentation in multiple languages, advertising the system's existence in different channels depending on where different groups tend to gather, or provide certain accessibility features.
2. User interface and user training. When designing, testing, and rolling out the system's user interface or pilot deployment, feedback should be gathered on its usability from every identify at-risk population group. Challenges raised by or about specific groups should be addressed in the development team's roadmap.
3. End user feedback. The AI system's development team should conduct studies to understand how end users perceive and interact with generated content and recommendations within the context of use. Since end users are not domain experts, and since expectations and actions are different within the content of use from those within a controlled testing environment, such studies can uncover insights and issues that do not come up in earlier stages. The team should proactively seek feedback from end users that reflect a wide range of demographics and backgrounds, specifically including those from

identified at-risk groups.

4. Automation bias. Automation bias is the human tendency to favor suggestions from automated decision-making systems and to ignore contradictory information. It is a human tendency to take the road of least cognitive effort – and research suggests that relying on human oversight to correct for biased AI algorithms does not work in practice. Our AI systems should monitor for automation bias. And where it exists should develop alternative models or user interface elements to mitigate it.
5. Unfair outcomes. In some cases, even after an AI system passed all fairness testing on its data, algorithms, and application, it may produce unfair outcomes due to unintended or unforeseen consequences. For example, a system which prioritized healthcare to sicker patients ended up being racially unfair, because it estimated which patients were sicker by how much they spent on healthcare, which favors people with good health insurance, which in turn is not equally distributed.

Therefore, our AI systems must monitor the actual outcomes produced for their stakeholders, for every at-risk population group. For example, regular reports should be able to answer, for each at-risk demographic group: How many patients received extra treatment? How many people got interviews, or got accepted to jobs? How many customer support issues were successfully resolved?

6. Drift. Model drift, also known as model decay, refers to the degradation of model performance over time because the data it was trained on doesn't match the real world anymore. Drift can result in faulty decision making and bad predictions. If not properly monitored over time, even the most well-trained, unbiased AI system can produce unwanted results in production.

Drift detection is a key component of AI governance. Ensuring fairness over time requires monitoring and correcting for drift in performance for each identified at-risk population, and not only the general validation dataset. Our systems should include monitoring for data and model drift for each at-risk population group, as well as procedures for correcting such drift by the AI system's development team.

## 9. Monitoring in Production

The system should implement real-time monitoring for analyzing how well the AI system serves different population groups. Monitors should enable measuring performance over time and across different user demographics, identify deviations from the desired standards, and trigger alerts for human intervention. The specific metrics to monitor depend on the type of the AI system and should aim to cover.

1. Access. Is there equal access to the system by different groups? For example, for a

system that matches candidate resumes to job listing, how many men versus women find about the system and upload their resume?

2. Usage. Are different groups able to use the system and complete key workflows in the same way? For example, how often do men and women complete the process to set up a job candidate resume, see matching jobs, and apply to them?
3. Decisions. Does the system recommend or decide similarly across groups? For example, do men and women get matched to open job listings at a similar rank?
4. Overrides. Do human users override the system's decisions, or edit the content it generated before using it, more often for some demographic groups?
5. Outcomes. Do people from different groups enjoy the same outcomes by using the system? For example, do men & women have similar success in finding a new job?
6. Drift. Is the proportion between demographic groups materially different from what it is in the AI system's training and validation datasets? Is it changing over time?

Depending on the AI system's Risk Level with respect to fairness, the monitoring system and metrics may be managed by the development team (Low), by a separate team within the organization (Managed), or by a qualified third party (Regulated).

# AI Transparency Policy

## 1. Purpose

This policy describes AI privacy controls in our organization (“we”, “us”, or “our”). It defines how we comply with relevant legislation and industry standards with respect to model transparency, explainability, data lineage, deepfakes, and required public disclosures.

## 2. Scope

This policy, in addition to applicable laws and regulations, is mandatory to us and our team members. We expect our business partners to meet the same standards and obligations as we demand from ourselves. This policy applies to all team members (whether temporary, fixed-term, or permanent) and third parties (contractors, seconded staff, trainees) acting under our responsibility. Together referred to as "team members". The policy also applies to our officers, trustees, and/or management at any level.

Any arrangements our company makes with a third party are subject to clear contractual terms, including specific provisions that require the third party to comply with minimum standards and procedures set out in this policy or other policies.

## 3. Disclose Use of AI

1. Any AI system that is intended to interact with people must always disclose to them that they are interacting with an AI system.
2. The user interface and underlying data model of an AI system must separate data generated by a human from data generated by AI. For example, if a medical system includes a field for ‘Diagnosis’, there must be two separate fields in the data model and user interface for a diagnosis entered by a human clinician versus a diagnosis predicted by the system.
3. If a system operates a bot – defined as “an automated online account where all or substantially all of the actions or posts of that account are not the result of a person” – then the bot must not communicate or interact with a person in order to incentivize a sale or transaction of goods or services, or to influence a vote in an election, without disclosing that it is a bot.
4. The AI system should inform users of its scope, precautions, and usage prohibitions in a user-friendly manner within contracts or service agreements, supporting informed choices and cautious use by users. A system’s API or user interface should include guardrails for user queries or data input that is outside the safe scope of the system. and respond with a relevant explanation of the issue.

5. In some high-risk applications, such as employment or criminal decisions, users may have the right to opt out of interacting with an AI. An AI system with such a requirement automatically has a Risk Level of Managed with respect to transparency. If the system supports this, it must make this option clearly visible and easily accessible to its users.
6. To comply with certain laws our organization may be required to make available a Model Documentation Form to competent authorities or downstream providers. The AI Governance Office must review and authorize any such document before it is made public or shared with authorities.

## 4. Disclose AI Generated Content

1. Clearly label AI generated content. People should be able to immediately recognize text, images, audio, video, or other content that was AI generated. This should be either made obvious in the user interface, or by adding the text “Disclaimer: this output has been generated by artificial intelligence” if the content can be consumed outside of the system’s user interface.
2. Watermarking. AI generated content intended to be published or consumed outside of a system’s internal user interface should be watermarked. Watermarking involves adding some detectable but not obvious mark to content, that lets third-party tools independently identify the provenance of AI generated content.

## 5. No Deepfakes

It is in some cases a criminal offense, and therefore at the Risk Level of Unacceptable, for an AI system to generate content that:

- Fabricates a deceptive video (“a real person performing an action that did not actually occur”) with intent to injure a candidate or influence the outcome of an election.
- Has images created, altered, adapted, or modified by electronic, mechanical, or other means to portray an identifiable minor engaged in sexual conduct.
- Fabricates intimate or sexually explicit images and depictions.
- Includes unauthorized creation and distribution of a person’s photograph, voice, or likeness.
- Includes commercial use of digital replicas of deceased performers in films, TV shows, video games, audiobooks, sound recordings, etc., without first obtaining the consent of those performers’ estates.

## 6. Explain AI Decisions

1. Explain AI decisions to users. An AI system should include the ability to explain why a decision was reached to its end users, in plain language appropriate to their level. For

example, if a person is denied a loan based on a model's decision, the system should generate an explanation of both the business logic and the input data points used in making that decision, and make that explanation easily accessible.

2. Define explainability needs. At the design phase, the requirements for explainability should be established with end users and domain experts. They should include (a) the goal of the explanations (e.g., global description of the model's behavior versus local explanation of each AI decision); (b) the most suitable approach for AI explainability; and (c) potential limitations to anticipate and monitor (e.g., over-reliance of end users on AI decisions).
3. Evaluate explainability. The explainable AI methods should be evaluated, first quantitatively by using computational methods to assess the correctness of the explanations, then qualitatively with end users to assess their impact on user satisfaction, confidence, and performance. The evaluations should also identify any limitations of the AI explanations: for example, are they incoherent, sensitive to noise or adversarial attacks, or unreasonably increase the confidence in the AI generated results?
4. Enable audits of AI decisions. An AI system should log the specific versions of the models, software, and data points used in making each individual decision. This should enable forensic analysis in case of audits, errors, or incidents that the system generated explanation is insufficient to resolve.

## 7. Disclosures When Acting as an AI Developer

1. Our organization is acting as an AI Developer when one of our teams is building an AI system designed to be used by other organizations. This includes either developing a new AI system or service or substantially modifying it, which means creating "a new version, new release, or other update to a generative artificial intelligence system or service that materially changes its functionality or performance, including the results of retraining or fine tuning."
2. An AI developer is required to document the training data used for its models:
  - Disclosure of where the training data is coming from, both for models trained or fine-tuned by the system and for underlying models such as foundation models the system is reusing.
  - Disclosure of any copyrighted materials included in the training data, and the legal basis or content license that allows its use for model training.
  - Disclosure of any personal information in the training data, including people's photograph, voice, or likeness, and the consent received to use this data for model training purposes.

3. An AI Developer should provide documentation describing the system’s intended use, known limitations, and material updates.
4. This documentation must be provided and reviewed by the AI Governance Officer as part of each system’s Pre-Deployment Checkpoint. An AI Developer of a system with a Risk Level of Managed or above may be required to provide a summary of the system’s Risk Assessment to its users. In such cases, the system’s Risk Level with respect to transparency is automatically Managed. This means that the system cannot be offered, sold, leased, given, or otherwise made available to a third party unless we provide it with sufficient information to perform a risk assessment on the use of the system, through a document detailing the potential risks, benefits, and intended uses of the system.
5. An AI Developer of a GPAI model is also required to provide:
  - A public and up-to-date Model Card or Model Documentation Form
  - Model documentation to downstream providers, and upon legitimate request to a competent national authority, provided that our organization is entitled to protect our intellectual property rights and confidential business information or trade secrets when providing such documentation.
  - Additional information within no later than 14 days upon a written request from downstream providers, when such information is necessary to enable them to have a good understanding of the capabilities and limitations of the GPAI model, relevant for its integration into the downstream providers’ AI system and to enable those downstream providers to comply with their regulatory obligations.
  - Maintain an up-to-date training data summary of an GPAI model.

Such documentation should be retained for ten years from the making the GPAI model available in the market. Before making a GPAI model available in the EU, it is also required to provide our organization’s contact information to the EU AI Office and to national authorities on request, as well as appoint an authorized representative in the EEA. The AI Governance Office must review and approve any such communications to regulatory authorities.

6. Model Cards. In some jurisdictions, laws require that a summary of a system’s Risk Assessment, training data sources, or performance and limitations must be made public online. This includes some models that make medical decisions, models that make employment decisions, GPAI models, GPAI models with systemic risk, and other cases as determined by the AI Governance Officer.

In such cases, this system’s Risk Level with respect to transparency is automatically Managed. The AI system’s team should collaborate with the AI Governance Officer on

defining the public model card, get its approval before making it public in whole or in part, and set procedures on how to update it when releasing new versions.

## 8. Disclosures When Acting as an AI Deployer

Our organization acts as an AI Deployer when it is licensing or reusing an AI system developed by another organization. We may still be required to document and disclose the same elements as when acting as an AI Developer – including a system’s Risk Assessment, source and licenses to training data, model performance and limitations.

Therefore, when acting as an AI Deployer, we must only deploy or reuse models from AI Developers who make available to us documentation:

1. Necessary to complete an impact assessment of the AI system
2. About the source and licenses to its training data
3. Summarizing the types of high-risk systems that the developer has developed and how reasonably foreseeable risks are managed
4. Disclosing reasonably foreseeable risks within 90 days after discovery of the risk
5. Describing intended uses, known limitations, and material updates.

As an AI Deployer, we shall submit a clear notice to an AI system’s end users when an automated decision-making technology is used to materially influence a consequential decision that impacts their safety, health, education, finances, or employment.

# AI Incident Reporting Policy

## 1. Purpose

This policy describes the processes for reporting AI incidents in our organization (“we”, “us”, or “our”). It defines how we identify, report, investigate and resolve AI incidents.

We aim to prevent AI incidents. If such an event, nevertheless, happens, we make sure to investigate the events and its causes, minimize its impact of on individuals and organizations, and take action to minimize the risk of its occurrence.

## 2. Scope

This policy, in addition to applicable laws and regulations, is mandatory for us and our team members. We expect our business partners to meet the same standards and obligations as we demand from ourselves. This policy applies to all team members (whether temporary, fixed-term, or permanent) and third parties (contractors, seconded staff, trainees) acting under our responsibility. The policy also applies to our officers, trustees, and/or management at any level.

Any arrangements our company makes with a third party are subject to clear contractual terms, including specific provisions that require the third party to comply with the minimum standards and procedures set out in this policy or other policies.

## 3. Internal Reporting

1. The AI Governance Officer should provide a central tool, which each AI system should have access to and utilize to record, respond to, and log AI incidents, including:
  - reports from individuals (e.g. as per the Whistleblower policy below) and end users;
  - reports from our team members, including AI developers, deployers, etc.;
  - reports from outside organizations.
2. An AI Incident should be reported to the AI Risk Manager as soon as possible, whereas our team members should report within 24 hours.
3. All information about each AI incident must be retained for a period of three years.
4. Each incident report should include, at a minimum, a System ID, Title, Reporter, System/Source, Data Reported, Date of Incident, Description, Impact(s), Impacted Stakeholders, Actions Taken, and Resolution Status. The description of each incident should also include:
  - Incident description and timeline;

- Harm details;
  - Impact on people and planet;
  - Economic context and impact;
  - Data and input dimension;
  - AI model dimension;
  - Task and output dimension;
  - Mitigation actions taken and planned, including owners and timeline.
5. Within 24 hours upon receipt of the report, the AI Risk Manager of the impact system will assign severity and priority level, assign an investigation and mitigation actions to our team members, define timeframes, and document all steps taken.
  6. The AI Risk Manager of each system is responsible for tracking that resolution, mitigation, remediation, and documentation of each incident is done.
  7. The AI Risk Manager should prepare a comprehensive resolution report within 60 days of the incident. Incident reports are intended for internal use to support audits, safety reviews, and continuous improvement. Special focus is given to whether the incident arose from system failure, malicious use, or unintentional misuse.
  8. The AI Risk Manager should implement corrective measures. Policy, process, and training gaps identified during the incident are addressed to prevent recurrence.
  9. The AI Risk Manager should notify relevant internal stakeholders, including executives, compliance and legal teams.
  10. All incident records must be securely retained for a minimum of **6 years**.

## 4. External Reporting

1. The AI Risk Manager should escalate certain AI incidents, and report to the Legal Counsel for further assessment.
2. The Legal Counsel will determine whether the AI incident must be reported to regulatory, supervisory, or enforcement bodies in any jurisdiction.
3. No information about an AI incident can be published or shared outside the organization, unless first reviewed and approved by Legal Counsel.
4. Upon assessment, the Legal Counsel should define whether an AI incident or other relevant event constitute:

- potential criminal offenses related to the use or impact of AI systems
  - cybersecurity threats
  - violation of data protection laws (e.g., GDPR in the EU, CCPA in California)
  - incidents posing threats to public health or safety
  - significant discriminatory or biased outcomes, particularly in regulated or high-impact domains such as healthcare, finance, hiring, or law enforcement
  - cases involving serious harm to individuals or communities, property, environment, critical infrastructure
5. Required Disclosures. Legal Counsel will determine if, how, to whom, and under what laws or regulations an incident must be externally disclosed. For example:
- The EU AI Act mandates that developers log any information about serious risks and accidents of their high-risk AI systems into a public database maintained by the EC.
  - Any serious incident identified in the course of the testing in real world conditions shall be reported to the EU national market surveillance authority. The provider or prospective provider shall adopt immediate mitigation measures or, failing that, shall suspend the testing in real world conditions until such mitigation takes place, or otherwise terminate it.
6. Voluntary Disclosures. Where there is no legal obligation to report externally, the AI Governance Officer may choose to report certain incidents to industry databases and repositories in the interest of public trust, ethical leadership, and contributing to collective learning within the AI community. Any information disclosed to this end must be first reviewed and approved by the AI Governance Officer and Legal Counsel, to ensure that it does not place undue legal, financial, or reputation risks to our organization.

## 5. Classification of Incidents

1. As define by the Center for Security and Emerging Technology’s (CSET), an “AI incident” is an event, circumstance or series of events where the development, use, or malfunction of one or more AI systems directly or indirectly leads to any of the following harms:
- injury or harm to the health of a person or groups of people;
  - disruption of the management and operation of critical infrastructure;
  - violations of human rights or a breach of obligations under the applicable law;
  - intended to protect fundamental, labor, and intellectual property rights; or
  - harm to property, communities, or the environment.

2. An “AI near miss” is defined as an event, circumstance, or series of events where the development, use, or malfunction of one or more AI systems could have directly or indirectly led to any of the following harms, but failed to by chance or was intercepted:
  - injury or harm to the health of a person or groups of people;
  - disruption of the management and operation of critical infrastructure;
  - violations of human rights or a breach of obligations under the applicable law;
  - intended to protect fundamental, labor, and intellectual property rights; or
  - harm to property, communities, or the environment.
3. An “AI hazard” is defined as an event, circumstance, or series of events where the development, use, or malfunction of one or more AI systems could plausibly lead to any of the following harms:
  - injury or harm to the health of a person or groups of people;
  - disruption of the management and operation of critical infrastructure;
  - violations to human rights or a breach of obligations under the applicable law
  - intended to protect fundamental, labor, and intellectual property rights; or
  - harm to property, communities, or the environment.
4. In this policy “AI incident” includes both incidents and near misses, except where a distinction is specifically noted. Both incidents and near misses should be reported, including cases of:
  - Algorithmic bias or discrimination;
  - Inaccurate or harmful outputs;
  - Privacy or data protection violations;
  - Unauthorized or unexpected system behavior;
  - Security vulnerabilities in AI models;
  - Failure to meet regulatory or ethical standards;
  - Misuse or abuse of AI systems

## 6. Whistleblower Policy

1. Our organization is committed to fostering a culture of transparency, accountability, and ethical conduct in the development and deployment of AI systems. To this end, we encourage employees and stakeholders to raise concerns regarding potential risks, misuse, or violations related to AI systems through internal reporting channels before considering

external escalation.

2. Internal Reporting Mechanisms. Employees may report concerns using:

- Anonymous reporting hotline;
- Reporting to our AI Risk Manager.

By establishing these mechanisms, we ensure individuals can share information safely without fear of retaliation. Reports through our hotline may be anonymous, protecting whistleblower identity.

3. Protection Against Retaliation. Retaliation against whistleblowers is prohibited. Employees who report reasonable concerns in good faith cannot be fired, demoted, harassed, or otherwise penalized due to making a report. Any retaliatory actions will themselves be treated as a policy violation subject to disciplinary measures.

4. Whistleblower Rewards. The AI Governance Office shall establish a reward program with financial incentives for individuals who disclose wrongdoing or AI-related risks leading to regulatory action.

# AI Copyright Policy

## 1. Purpose

This policy establishes principles, processes, and controls in our organization (“Company”, “we”, “us”, or “our”) to ensure that all AI systems developed, procured, or deployed by the organization comply with and respect copyright law.

## 2. Scope

This policy applies to any:

- AI System
- GPAI model
- open-source AI models, systems, or GPAI models
- High-Risk AI Systems or GPAI models with systemic risk

This policy applies to all AI developers, deployers, and downstream providers. Each covered party shall:

- Comply with applicable AI copyright laws and regulations,
- Adhere to obligations related to the provisions, deployment, modification, and/or distribution of covered AI products.

## 3. Web Crawling Controls

- Our obligations related to deploying web crawling systems:
- Web crawling systems deployed for training purposes must limit access to legally available content only.
- Circumventing technical barriers (such as paywalls or access restriction mechanisms) is strictly prohibited.
- Exclude from crawling websites that US, EU, and EEA authorities have recognized as persistently and repeatedly infringing copyright and related rights on a commercial scale, based on a publicly maintained registry.
- Comply with machine-readable rights reservations by copyright holders. This includes employing web-crawlers that follow robots.txt and follow other widely adopted machine-readable protocols.
- Commit to transparency regarding crawler obligations, our approach to handling rights-reserved content, with automated notifications for rights holders.

## 4. Mitigating the Risk of Infringing Outputs

- AI Systems shall implement appropriate and proportionate technical safeguards to minimize the likelihood that AI models produce copyright-infringing content.
- We prohibit infringing uses, including for open-source models, whenever we use an AI model directly or through third-party implementation.
- Any GPAI models we provide shall be accompanied by documentation alerting users to the prohibition of copyright infringing uses of the model.

## 5. Complaint Mechanism

Copyright holders may submit complaints regarding alleged copyright infringement connected to our AI systems, by sending a detailed email to [legal@pacific.ai](mailto:legal@pacific.ai).

We will acknowledge complaints promptly. Substantive responses and any corrective measures will be provided within a reasonable timeframe, in accordance with applicable law and contractual obligations.

Where necessary, we will suspend or adjust AI system functionality pending investigation.

# AI Acceptable Use Policy

## 1. Purpose

The purpose of this AI Acceptable Use Policy is to define the conditions under which artificial intelligence systems may be developed, acquired, deployed, and used. The acceptable use of AI must align with our mission and values, support the welfare and rights of our stakeholders and the public, and comply with all applicable laws and regulations.

All AI use should contribute positively to society. Uses of AI that contravene these objectives, violate legal obligations, or present unreasonable risks to individuals, communities, or the environment are strictly prohibited.

## 2. Scope

This policy is mandatory for all internal team members, business partners, and customers who access or deploy AI systems that we develop, deploy, or license.

- **Team Members:** This includes all employees, contractors, seconded staff, interns, trainees, and other individuals performing work under the organization's direct authority—whether on a permanent, fixed-term, or temporary basis.
- **Business Partners:** This includes vendors, consultants, service providers, and any third parties acting on behalf of the organization or providing AI-related services. Any contractual agreement with such parties should include provisions requiring full compliance with this policy and its underlying procedures.
- **Customers:** Where this policy is referenced within a license agreement, service contract, or End User License Agreement (EULA), it becomes binding on customers who use our AI systems or outputs. Customers are expected to adhere to the same standards of responsible AI use and ethical deployment.

Together, these groups are referred to collectively as "users" under this policy framework.

## 3. Unacceptable Uses of AI

Each AI system must undergo a formal risk classification. Risk mitigation, testing, and required approvals are determined by the assigned risk level. No system classified as "unacceptable" may be developed or deployed. This risk-based approach ensures that AI systems are designed, operated, and governed in proportion to their potential harms and societal impact.

The following categories of AI applications are classified as unacceptable and are strictly prohibited under this policy. This list unifies unaccepted uses as defined by US federal and state laws with the acceptable use policies of major cloud and model providers.

### 3.1 Human Rights, Civil Liberties, and Safety

- Autonomous weapons or lethal autonomous systems
- Predictive policing or pre-emptive criminal risk scoring based solely on profiling or assessing the relevant individual's personality traits and characteristics
- Social scoring systems that rate individuals based on behavior, socioeconomic status, or characteristics
- Biometric categorization or real-time biometric identification in public spaces
- Unlawful tracking or stalking systems
- Surveillance systems used to target or discriminate against, harass a person or vulnerable populations
- Systems for ongoing surveillance or real-time or near real-time identification or persistent tracking of the individual using any of their personal data, including biometric data, without the individual's express consent
- Voice-activated or generative systems that manipulate children's behavior in harmful ways
- Advertising using AI generated actors or voice without disclosure
- Use of AI to facilitate human trafficking or the exploitation of vulnerable individuals
- Training or deploying AI to engage in or assist fraud, extortion, or identity theft
- Unauthorized practice of any regulated profession, such as legal, financial, therapeutic or medical profession
- Inferring the emotions of a natural person in the workplace and educational institutions, except for medical or safety reasons

### 3.2 Misinformation, Influence, and Deception

- AI systems used to manipulate, deceive, or impersonate others without consent
- Systems that produce or disseminate misinformation, disinformation, or propaganda
- Deepfakes or synthetic media used without consent or disclosure, or for harassment, sexual exploitation, fraud, or election interference
- Deepfakes or sexual scenes with a minor that abuse or harass the minor
- Bots that interact with humans for commercial or political purposes without disclosing their automated nature
- Use of AI to manipulate elections, influence voter behavior, or suppress participation
- AI-generated content that encourages or instructs on self-harm or suicide

- Any content that incites violence or physical harm against others
- Use of AI to create or distribute content that promotes terrorism or violent extremism
- Use of AI for the creation or dissemination of intimate imagery without consent (non-consensual deepnudes)
- Use of AI to generate fake product or service reviews
- Use of AI for astroturfing campaigns (false grassroots support)
- Training or deploying AI tools that impersonate emergency services or medical providers

### 3.3 Data Privacy, Consent, and Security

- Re-identification of anonymized data without consent
- Use of personal data for training AI without informed, specific, and revocable consent
- Training on datasets that include copyrighted, confidential, or sensitive material without proper licenses or legal basis
- AI systems that violate privacy or data protection laws, including unauthorized surveillance, profiling, or scraping
- Generative models that memorize and reproduce sensitive or proprietary data
- Facial recognition databases without express consent

### 3.4 Discrimination and Unfair Outcomes

- AI systems that produce or reinforce discriminatory outcomes based on race, gender, ethnicity, religion, age, disability, or other protected characteristics
- Employment, financial, medical, or educational decision-making systems without bias evaluation and audit trails
- Systems that enable digital redlining or deny opportunities unfairly to protected groups
- Use of demographic attributes for personalization without regulatory justification and independent review

### 3.5 Intellectual Property and Ethical Content Generation

- Generation of copyrighted material without licensing or legal basis
- Creation of AI-generated actors, voices, or likenesses without disclosure or consent
- Commercial use of deceased individuals' likenesses without estate permission
- Fabrication of scientific research, legal documents, or medical advice
- Generation of content that is obscene, offensive, harmful, or unsafe for minors

### 3.6 Safety and Misuse Prevention

- Deployment of AI systems that present unmitigated risks to critical infrastructure, health, or public safety
- Use of AI in life-critical settings without extensive validation, monitoring, and fail-safes
- Circumvention of safety measures, filters, or integrity mechanisms in AI models or APIs
- Use of AI for spam, phishing, malware generation, or security evasion
- Training AI systems to generate exploits, jailbreaks, or instructions for circumventing security software
- Using AI to assist in denial-of-service attacks or intrusion into protected networks
- Using AI to impersonate licensed professionals without explicit user awareness and disclaimers

## 4. Enforcement and Review

Any AI system identified as violating this Acceptable Use Policy will be immediately subject to a mandatory stop. The system must undergo a risk reassessment by the AI Governance Officer. If the use is determined to fall under an unacceptable category, it must be terminated, and any deployment reversed.

Please submit any questions, feedback, or concerns to our team at [legal@pacific.ai](mailto:legal@pacific.ai)



Pacific AI

[www.pacific.ai](http://www.pacific.ai)

[info@pacific.ai](mailto:info@pacific.ai)