

A futuristic, metallic robot with glowing blue eyes and a digital interface in the background. The robot is standing in a dark, futuristic environment with a glowing blue and green digital interface on the left. The interface displays a complex circuit board pattern and a grid of data points. The robot's body is composed of various mechanical parts and panels, with a glowing blue light emanating from its chest area. The overall scene is set against a dark, atmospheric background with some light flares and a faint, repeating pattern of the robot's head in the distance.

AI Governance and Regulation: Aspects of Ethics, Privacy, and Security

February 13, 2024
ISACA L.A. Chapter

Katharina Koerner, PhD
Responsible AI Advisor - Tech Diplomacy Network
Corporate Development Manager - Daiki



Presentation Overview

1. Responsible AI principles.
2. Existing and upcoming AI regulation.
3. AI management and governance frameworks.
4. AI security.
5. Privacy challenges and solutions.

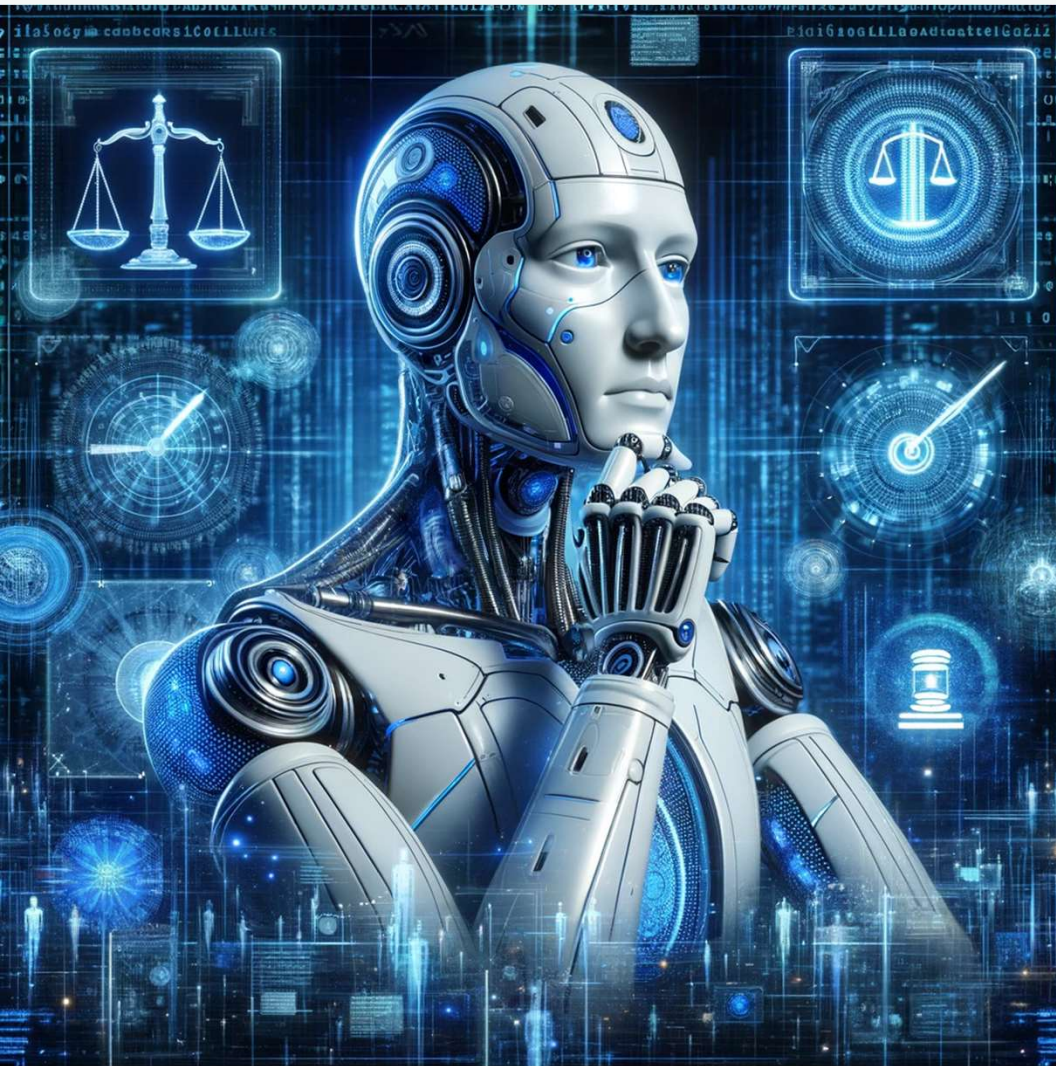
1. Responsible AI Principles

The terms “Ethical AI”, “Trustworthy AI” and “Responsible AI” are often used interchangeably. For others, ethics goes beyond or is different from RAI.

Responsible AI frameworks are developed and implemented as self-regulatory initiatives, by international organizations, and standardization bodies.

Existing AI regulation can regularly be mapped to the principles of responsible AI.





2. AI Regulation

US:

- AI regulated in sectoral approach, e.g., FTC, EEOC, CFPB, ..
- Executive Order on AI (not as regulation per se)
- State (privacy) laws

EU + extraterritorial scope:

- GDPR for personal data
- EU AI Act (upcoming)
- EU Liability directive (upcoming)

Other countries:

- Canada, China, Brazil,...

Global AI Governance initiatives:

- UN, G7, OECD,...

ENFORCEMENT EXAMPLES - US

The FTC's biggest AI enforcement tool? Forcing companies to delete their algorithms

Algorithm disgorgement requires companies to remove products built on data they shouldn't have used in the first place.

BY TONYA RILEY • JULY 5, 2023



News

EEOC settles first AI-bias case, but many more are likely

Employment attorneys caution that companies will be on the hook if a vendor's hiring software turns out to perpetuate bias. While such systems can be complicated and difficult to understand, ignorance is no excuse.

By Greg Andrews | August 16, 2023 at 11:13 AM



AP



U.S.

WORLD

POLITICS

VIDEO

SPOTLIGHT

SPORTS

BUSINESS

SCIENCE

FACT CHECK



Israel-Hamas war

Brussels shooting

Trump gag order

Colorado train der

FTC investigating ChatGPT creator OpenAI over consumer protection issues



Generative AI refers to a class of artificial intelligence (AI) models that can create or generate new data, such as images, text, or music, that is similar to the data it was trained on. It identifies patterns and relationships in the input data and then

ENFORCEMENT EXAMPLES - EU



Join TechCrunch+

Login

Search Q

TechCrunch+

Startups

Venture

Security

AI



REUTERS®

World ▾

Business ▾

Markets ▾

Sust

Uber still dragging its feet on algorithmic transparency, Dutch court finds

Natasha Lomas @riptari / 11:00 AM PDT • October 5, 2023



OpenAI's ChatGPT breaches privacy rules, says Italian watchdog

Reuters
January 30, 2024 12:09 AM PST · Updated 15 days ago

ARTIFICIAL INTELLIGENCE

ChatGPT



Boards, Policy & Regulation | Data Privacy | Litigation

Poland investigates OpenAI over privacy concerns

Reuters

September 21, 2023 2:10 PM PDT · Updated a month ago

August 13, 2021

Italian data protection supervisory authority fines two food delivery companies for non-compliant algorithmic processing

in LinkedIn

f Facebook

t Twitter

Send

Embed

UPCOMING EU AI ACT: SCOPE

Compliance with the EU AI Act is mandatory when targeting or impacting the EU market, regardless of the company's location.

Specifically applicable globally for:

Providers (Developers or Marketers of AI Systems) when

- AI systems are introduced to the EU market.
- AI system's output is used in the EU under international law.

Deployers (Users of AI Systems) when

- International law applies or the system's output is used in the EU.

UPCOMING EU AI ACT: requirement overview

Prohibited Practices / Unacceptable Risk::

- Indiscriminate biometric data scraping.
- Public facial and emotion recognition systems.
- Social scoring by authorities.
- Behavior-manipulating AI.
- Profiling for delinquency.

High Risk:

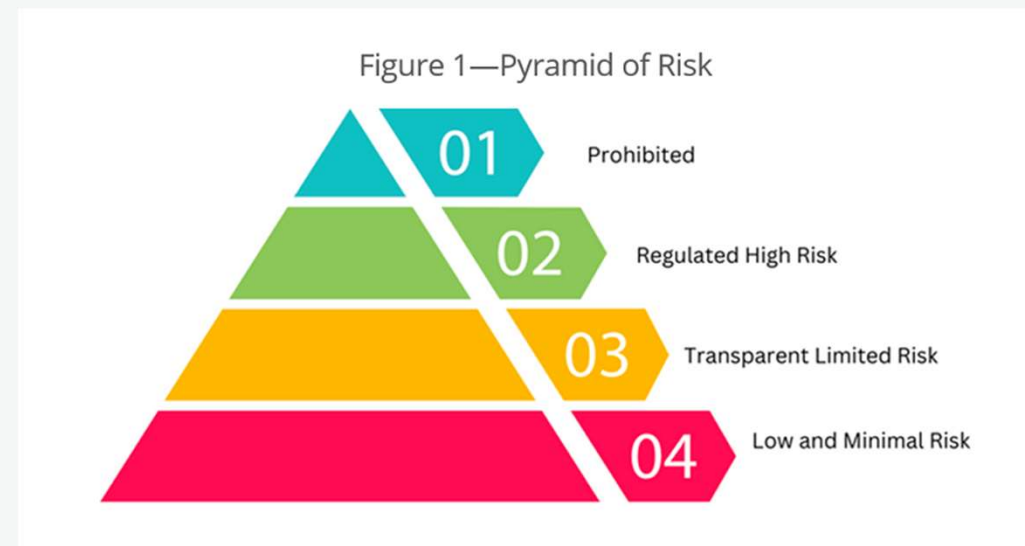
- Affects human safety and fundamental rights.
- Strict vetting and continuous monitoring.
- Mandatory EU database registration.

Limited Risk:

- Transparency about AI use and data type.

Minimal Risk:

- Encourages ethical use codes.



Source: ISACA

RISKS



SOURCE: PwC

EXAMPLES - REGULATORY RISK

Grading Foundation Model Providers' Compliance with the Draft EU AI Act











Source: Stanford Center for Research on Foundation Models (CRFM), Institute for Human-Centered Artificial Intelligence (HAI)

	OpenAI	cohere	stability.ai	ANTHROPIC	Google	BigScience	Meta	AI21labs	ALEPH ALPHA	ELEUTHERAI	Totals
Draft AI Act Requirements	GPT-4	Cohere Command	Stable Diffusion v2	Claude 1	PaLM 2	BLOOM	LLaMA	Jurassic-2	Luminous	GPT-NeoX	
Data sources	● ○ ○ ○ ○	● ● ● ● ○	● ● ● ● ●	○ ○ ○ ○ ○	● ● ● ○ ○	● ● ● ● ●	● ● ● ● ●	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ● ●	22
Data governance	● ● ● ○ ○	● ● ● ○ ○	● ● ● ○ ○	○ ○ ○ ○ ○	● ● ● ○ ○	● ● ● ● ●	● ● ● ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ○ ○	19
Copyrighted data	○ ○ ○ ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ● ●	7
Compute	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ● ●	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ● ●	● ● ● ● ●	○ ○ ○ ○ ○	● ○ ○ ○ ○	● ● ● ● ●	17
Energy	○ ○ ○ ○ ○	● ○ ○ ○ ○	● ● ● ● ●	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ● ●	● ● ● ● ●	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ● ●	16
Capabilities & limitations	● ● ● ● ●	● ● ● ● ○	● ● ● ● ●	● ○ ○ ○ ○	● ● ● ● ●	● ● ● ● ○	● ● ● ○ ○	● ● ● ○ ○	● ○ ○ ○ ○	● ● ● ● ○	27
Risks & mitigations	● ● ● ● ○	● ● ● ○ ○	● ● ● ○ ○	● ○ ○ ○ ○	● ● ● ○ ○	● ● ● ○ ○	● ○ ○ ○ ○	● ● ● ○ ○	○ ○ ○ ○ ○	● ○ ○ ○ ○	16
Evaluations	● ● ● ● ●	● ● ● ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ○ ○	● ● ● ○ ○	● ● ● ○ ○	○ ○ ○ ○ ○	● ○ ○ ○ ○	● ○ ○ ○ ○	15
Testing	● ● ● ● ○	● ● ● ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ○ ○	● ● ● ○ ○	○ ○ ○ ○ ○	● ○ ○ ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	10
Machine-generated content	● ● ● ● ○	● ● ● ○ ○	○ ○ ○ ○ ○	● ● ● ○ ○	● ● ● ○ ○	● ● ● ● ●	○ ○ ○ ○ ○	● ● ● ○ ○	● ● ● ○ ○	● ● ● ○ ○	21
Member states	● ● ● ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ○ ○	● ● ● ● ●	○ ○ ○ ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ○ ○ ○ ○	● ● ● ○ ○	9
Downstream documentation	● ● ● ● ○	● ● ● ● ●	● ● ● ● ●	○ ○ ○ ○ ○	● ● ● ● ●	● ● ● ● ●	● ● ● ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ○ ○	24
Totals	25 / 48	23 / 48	22 / 48	7 / 48	27 / 48	36 / 48	21 / 48	8 / 48	5 / 48	29 / 48	

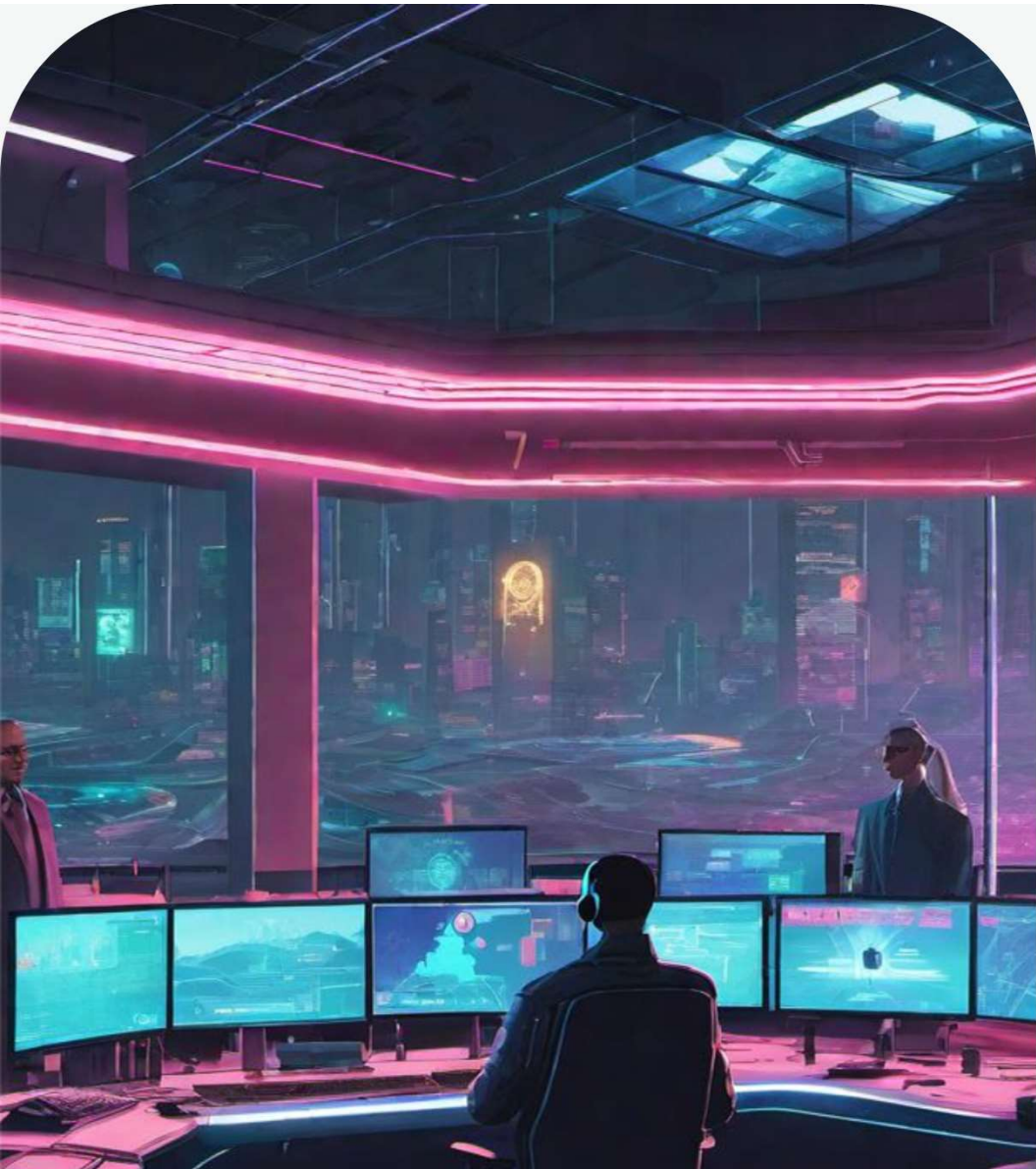
EXAMPLES - REGULATORY RISK

Foundation Model Transparency Index Scores by Major Dimensions of Transparency, 2023

Source: 2023 Foundation Model Transparency Index

	 Meta	 BigScience	 OpenAI	 stability.ai	 Google	 ANTHROPIC	 cohere	 AI21labs	 Inflection	 amazon	Average
	Llama 2	BLOOMZ	GPT-4	Stable Diffusion 2	PaLM 2	Claude 2	Command	Jurassic-2	Inflection-1	Titan Text	
Data	40%	60%	20%	40%	20%	0%	20%	0%	0%	0%	20%
Labor	29%	86%	14%	14%	0%	29%	0%	0%	0%	0%	17%
Compute	57%	14%	14%	57%	14%	0%	14%	0%	0%	0%	17%
Methods	75%	100%	50%	100%	75%	75%	0%	0%	0%	0%	48%
Model Basics	100%	100%	50%	83%	67%	67%	50%	33%	50%	33%	63%
Model Access	100%	100%	67%	100%	33%	33%	67%	33%	0%	33%	57%
Capabilities	60%	80%	100%	40%	80%	80%	60%	60%	40%	20%	62%
Risks	57%	0%	57%	14%	29%	29%	29%	29%	0%	0%	24%
Mitigations	60%	0%	60%	0%	40%	40%	20%	0%	20%	20%	26%
Distribution	71%	71%	57%	71%	71%	57%	57%	43%	43%	43%	59%
Usage Policy	40%	20%	80%	40%	60%	60%	40%	20%	60%	20%	44%
Feedback	33%	33%	33%	33%	33%	33%	33%	33%	33%	0%	30%
Impact	14%	14%	14%	14%	14%	0%	14%	14%	14%	0%	11%
Average	57%	52%	47%	47%	41%	39%	31%	20%	20%	13%	

Scores for 10 major foundation model developers across 13 major dimensions of transparency.



AI MANAGEMENT FRAMEWORKS

... help organizations navigate their involvement with AI systems, set up governance processes, and manage AI risks.

... streamline processes for using, developing, monitoring, or providing AI-related products and services.

... address trustworthiness concerns (security, safety, fairness, transparency, data quality).

... provide guidelines for deploying controls to support AI management processes.

... generate evidence of responsibility and accountability regarding the organization's role with AI systems.



AI GOVERNANCE

... is a component of AI management.

... focuses on establishing policies, guidelines, and ethical considerations for the responsible development, deployment, and use of AI within an organization.

... emphasizes the need for ethical and responsible practices in the use of AI technologies (e.g., in defined AI Objectives).

ISO/IEC 42001:2023, Artificial intelligence Management system

ISO/IEC 42001 is a certifiable framework that helps to establish, implement, maintain, continually improve and document an AIMS.

Goal:

Guidance for organizations of any size or type in the responsible development, provision, or utilization of AI systems to meet organizational objectives, applicable requirements, and obligations and expectations related to stakeholders.

Contents		Page
	Foreword.....	v
	Introduction.....	vi
1	Scope.....	1
2	Normative references.....	1
3	Terms and definitions.....	1
4	Context of the organization.....	5
4.1	Understanding the organization and its context.....	5
4.2	Understanding the needs and expectations of interested parties.....	6
4.3	Determining the scope of the AI management system.....	6
4.4	AI management system.....	6
5	Leadership.....	7
5.1	Leadership and commitment.....	7
5.2	AI policy.....	7
5.3	Roles, responsibilities and authorities.....	8
6	Planning.....	8
6.1	Actions to address risks and opportunities.....	8
6.1.1	General.....	8
6.1.2	AI risk assessment.....	9
6.1.3	AI risk treatment.....	9
6.1.4	AI system impact assessment.....	10
6.2	AI objectives and planning to achieve them.....	10
6.3	Planning of changes.....	11
7	Support.....	11
7.1	Resources.....	11
7.2	Competence.....	11
7.3	Awareness.....	12
7.4	Communication.....	12
7.5	Documented information.....	12
7.5.1	General.....	12
7.5.2	Creating and updating documented information.....	12
7.5.3	Control of documented information.....	13
8	Operation.....	13
8.1	Operational planning and control.....	13
8.2	AI risk assessment.....	13
8.3	AI risk treatment.....	14
8.4	AI system impact assessment.....	14
9	Performance evaluation.....	14
9.1	Monitoring, measurement, analysis and evaluation.....	14
9.2	Internal audit.....	14
9.2.1	General.....	14
9.2.2	Internal audit programme.....	14
9.3	Management review.....	15
9.3.1	General.....	15
9.3.2	Management review inputs.....	15
9.3.3	Management review results.....	15
10	Improvement.....	15
10.1	Continual improvement.....	15
10.2	Nonconformity and corrective action.....	16
	Annex A (normative) Reference control objectives and controls.....	17

AI Risk Management Framework



NIST AI RISK MANAGEMENT FRAMEWORK



MAP 1.5

Organizational risk tolerances are determined and documented.

About

Risk tolerance reflects the level and type of risk the organization is willing to accept while conducting its mission and carrying out its strategy.

Organizations can follow existing regulations and guidelines for risk criteria, tolerance and response established by organizational, domain, discipline, sector, or professional requirements. Some sectors or industries may have established definitions of harm or may have established documentation, reporting, and disclosure requirements.

Within sectors, risk management may depend on existing guidelines for specific applications and use case settings. Where established guidelines do not exist, organizations will want to define reasonable risk tolerance in consideration of different sources of risk (e.g., financial, operational, safety and wellbeing, business, reputational, and model risks) and different levels of risk (e.g., from negligible to critical).

Risk tolerances inform and support decisions about whether to continue with development or deployment – termed “go/no-go”. Go/no-go decisions related to AI system risks can take stakeholder feedback into account but remain independent from stakeholders’ vested financial or reputational interests.

If mapping risk is prohibitively difficult, a “no-go” decision may be considered for the specific system.

Suggested Actions

- Utilize existing regulations and guidelines for risk criteria, tolerance and response established by organizational, domain, discipline, sector, or professional requirements.
- Establish risk tolerance levels for AI systems and allocate the appropriate oversight resources to each level.

- Establish risk criteria in consideration of different sources of risk, (e.g., financial, operational, safety and wellbeing, business, reputational, and model risks) and different levels of risk (e.g., from negligible to critical).
- Identify maximum allowable risk tolerance above which the system will not be deployed, or will need to be prematurely decommissioned, within the contextual or application setting.
- Articulate and analyze tradeoffs across trustworthiness characteristics as relevant to proposed context of use. When tradeoffs arise, document them and plan for traceable actions (e.g.: impact mitigation, removal of system from development or use) to inform management decisions.
- Review uses of AI systems for “off-label” purposes, especially in settings that organizations have deemed as high-risk. Document decisions, risk-related trade-offs, and system limitations.

Transparency & Documentation

Organizations can document the following:

- Which existing regulations and guidelines apply, and the entity has followed, in the development of system risk tolerances?
- What criteria and assumptions has the entity utilized when developing system risk tolerances?
- How has the entity identified maximum allowable risk tolerance?
- What conditions and purposes are considered “off-label” for system use?

AI Transparency Resources

- [GAO-21-519SP: AI Accountability Framework for Federal Agencies & Other Entities.](#)
- [WEF Model AI Governance Framework Assessment 2020.](#)
- [WEF Companion to the Model AI Governance Framework- 2020.](#)

PRACTICAL IMPLEMENTATION

References

[Board of Governors of the Federal Reserve System. SR 11-7: Guidance on Model Risk Management. \(April 4, 2011\).](#)

[The Office of the Comptroller of the Currency. Enterprise Risk Appetite Statement. \(Nov. 20, 2019\).](#)

[Brenda Boulwood. How to Develop an Enterprise Risk-Rating Approach \(Aug. 26, 2021\). Global Association of Risk Professionals \(garp.org\). Accessed Jan. 4, 2023.](#)

Virginia Eubanks. 1972. Automating Inequality: How High-Tech Tools Profile, Police, and Punish the Poor. New York, NY, St. Martin's Press, 2018.

[GAO-17-63: Enterprise Risk Management: Selected Agencies' Experiences Illustrate Good Practices in Managing Risk. \(See Table 3\).](#)

[NIST Risk Management Framework.](#)

Using the AI Risk Management Framework

Workday

The voluntary NIST AI Risk Management Framework was developed through a collaborative process by industry, the society, academic, and government stakeholders. The framework is designed to equip organizations and individuals with approaches that increase the trustworthiness of AI systems, and to help foster their responsible design, development, and deployment. NIST does not validate or endorse any individual organization or its approach to using the AI RMF.

Benefits of Using the Framework

Workday is using the AI RMF to assess, refine, and strengthen our approach to trustworthy AI. Workday has benchmarked our common control framework to the AI RMF, augmenting guidelines, policies, and procedures and developing plans for further alignment where needed.

Workday has anchored our new responsible AI guidelines, product risk evaluation, and third-party risk questionnaire in the AI RMF.

Workday has used the AI RMF as a common reference point for Product & Technology, Responsible AI, Product Legal, and Data Privacy & Engineering teams to collaborate on AI risk management.

Situation & Drivers

Workday is a provider of enterprise cloud applications and an AI developer. Our applications for financial management, human resources, planning, spend management, and analytics are built with AI and ML at the core to help organizations around the world embrace the future of work.

Workday is used by more than 10,000 organizations around the world and across industries—from medium-sized businesses to more than 50% of the Fortune 500.

Our customers trust us with some of their most sensitive data, and upholding that trust is essential to Workday's business. Earning and retaining our customers' trust is aligned with our core values and facilitates the adoption of innovative AI tools by the world's leading organizations.

When Workday uses an AI tool that we have either developed or procured from a third-

party and interacts with end-users, we a deployer.

With AI best practices and technical star continuing to mature, the AI RMF is an benchmark for AI developers and deploy Workday uses the AI RMF to assess and our responsible AI practices and to come the rigor of our AI governance approach enterprise customers.

With leaders in the U.S., Europe, Canada, Singapore, and elsewhere taking action promote trustworthy AI, the AI RMF is a language to understand the expectations requirements of AI governance.

"As business leaders begin their AI journey, many are looking for a roadmap for how to develop and use AI in a way that is responsible and innovative. At Workday, we've found NIST's AI Risk Management Framework to be a concrete benchmark for mapping, measuring, and managing our approach to AI governance. We believe the Framework will help us maintain our customers' trust and stay our company's more resilient as we leverage innovation going forward."

Jim Stratton, Chief Technology Officer,

Process

Workday was an early champion of the AI and participated in every stage of its open multi-stakeholder development process. NIST launched the AI RMF 1.0 in January. Workday's Co-President Sayan Chakrabarti endorsed it as a "major milestone" in AI governance.

With support from senior leadership, and informed by our experience using the NIST Cybersecurity and Privacy Frameworks, we began using the AI RMF to map, measure, manage, and govern potential AI risk.

Workday's Privacy & Data Engineering team mapped the AI RMF to our existing common control frameworks. In doing so, they identified existing controls and processes correspond to the AI RMF's categories and subcategories.

Machine Learning

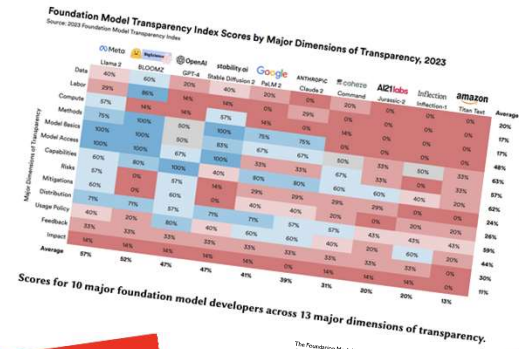
Fairness: Types of Bias

Estimated Time: 5 minutes

Machine learning models are not inherently objective. Engineers train models by feeding them a data set of training examples, and human involvement in the provision and curation of this data can make a model's predictions susceptible to bias.

When building models, it's important to be aware of common human biases that can manifest in your data, so you can take proactive steps to mitigate their effects.

WARNING: The following inventory of biases provides just a small selection of biases that are often uncovered in machine learning data sets; this list is not intended to be exhaustive. [1 of countless biases](#) enumerates over 100 different types of human bias that should be on the lookout for any and all potential sources of



Transparency and explainability of AI systems from ethical guidelines to requirements. *Nagendra Babu, Arun Kumar, Manoj Kumar, Anil Kumar, Karthikeyan, Sanjiv Kumar, July 2022*

Busting Companies for AI Bias in Hiring Is Tough Task for EEOC

DEEP DIVE

Riddhi Setty, Reporter

Annelise Gilbert

- Difficulty in identifying AI bias could inhibit agency
- Legislation offers possible avenue to needed information

The EEOC's high-profile effort to target biased hiring decisions made by companies' artificial intelligence software is poised to face big challenges, as the opaque nature of these tools often prevents applicants from knowing they're the victims of discrimination at all.

The US Equal Employment Opportunity Commission settled its first-ever AI-based hiring discrimination case on Thursday. A group of job seekers received \$365,000 in the deal, resolving the commission's suit against TutorGroup Inc., a company that allegedly programmed its recruitment software to automatically reject other applicants.

The Register

Funnily enough, AI models must follow privacy law – including right to be forgotten

We'll just take a look at the training data, oh... wait

Thu 13 Jul 2023 01:17 UTC

Thomas Claburn

Right to be Forgotten in the Era of Large Language Models: Implications, Challenges, and Solutions

Dawen Zhang^{1,2}, Pamela Finckenberg-Broman¹, Thong Hoang¹, Shidong Pan^{1,2}, Zhenchang Xing^{1,2}, Mark Staples¹, Xiwei Xu¹

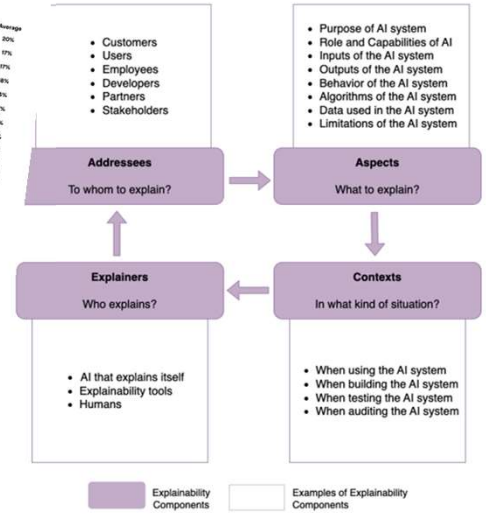


Fig. 2. A model of explainability components.

OPERATIONAL CHALLENGES

FY 2022

Ethical AI Startup Landscape



Disclaimer: Logos sourced from search engines or company websites. EAIDS was not able to obtain consent from every firm represented in this graphic, but is actively working to do so. Logos are presented here in good faith as a positive image to each company, but can be removed upon request.

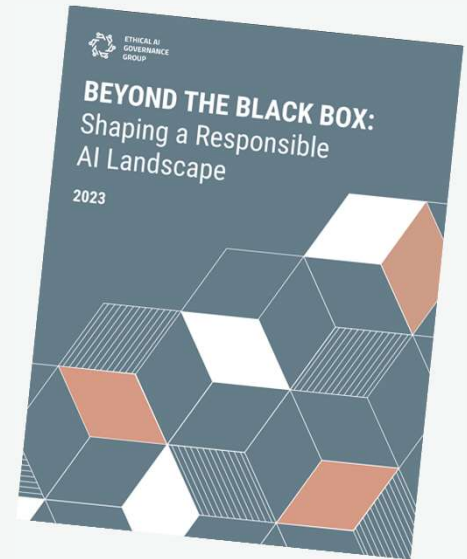
Don't see your company here? Make a [submission!](#)

AN EMERGING ECOSYSTEM AROUND RESPONSIBLE AI

The landscape of AI startups providing **services for responsible AI** operationalization is growing fast.

Includes:

- Identifying AI Deployment Across Organizations
- Ensuring Compliance with Regulatory Requirements during development
- Privacy-first data and processing



4. AI Security

Field is in development.

New resources:

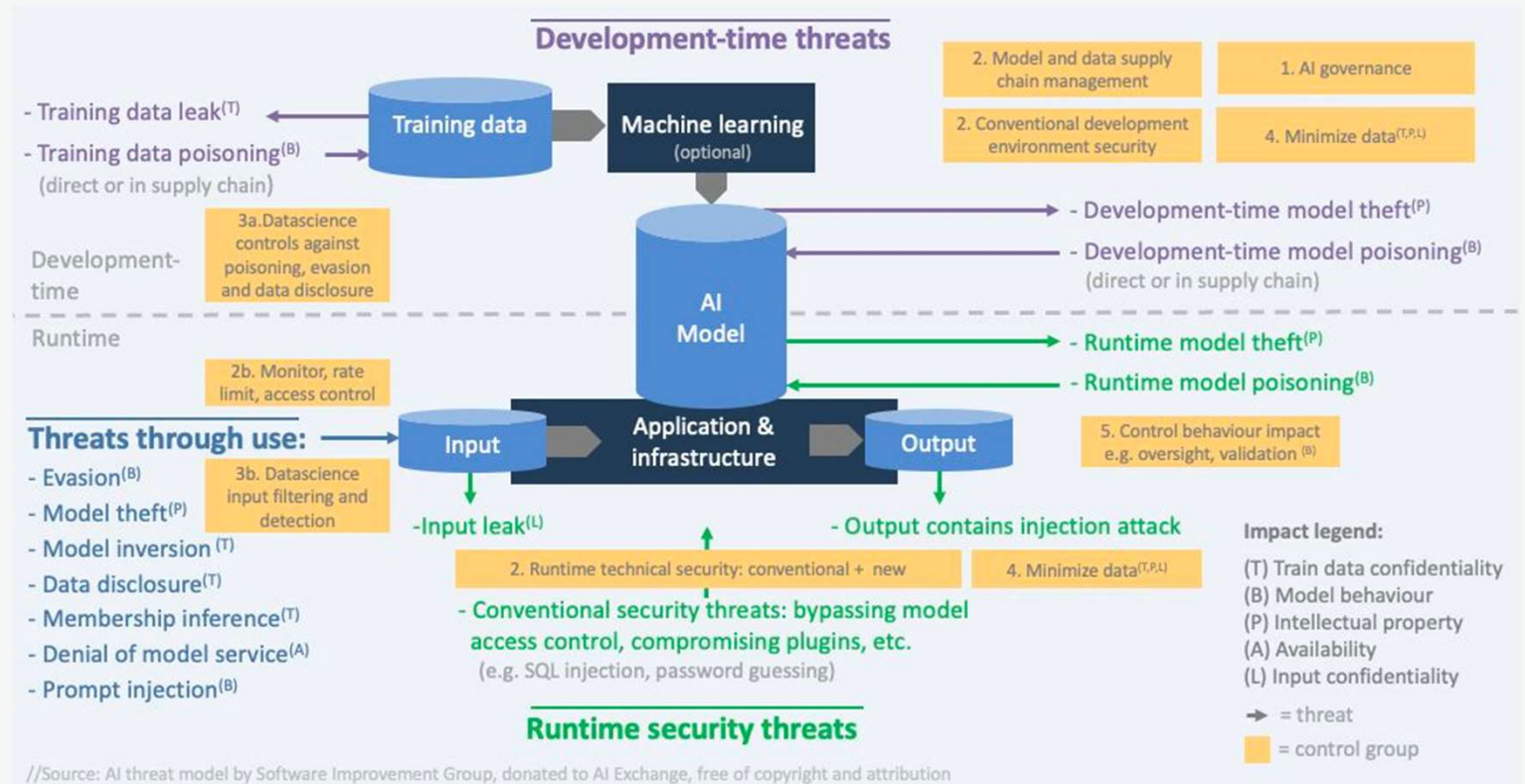
NIST Adversarial Machine Learning: A Taxonomy and Terminology of Attacks and Mitigations

ENISA Securing Machine Learning Algorithms

Microsoft Failure Modes in Machine Learning

Google Secure AI Framework (SAIF)

The OWASP AI Exchange



Essentials of AI Security

1. Well-established information security management systems form the foundation of AI security.
1. Create explicit and supplementary AI security measures and incorporate them into current risk management protocols.
1. Effective communication between cybersecurity experts and data scientists is essential for success in this arena.



5. Privacy Challenges in AI

General Privacy and Data protection principles apply to AI systems and processing of personal data.

AI privacy risks include:

- *Bias and Discrimination*
- *Harms from Inferences*
- *Problematic Data Actions: (unexpected data collection, storage, and use)*
- *Lack of Processing Basis: Processing data without a clear legal basis*
- *Secondary Data Use: Data repurposed for new, undisclosed uses*
- *Jurisdictional Issues and Data Scraping: Data scraping across jurisdictions*
- *Re-identification Risks in Anonymized Data*

Effective management of these risks requires adherence to privacy principles, clear regulatory frameworks, mechanisms for accountability, and engineering privacy into AI system development.



PRIVACY ENHANCING TECHNOLOGIES IN AI/ML

PETs that support privacy and security in the context of AI/ML and responsible AI are also referred to as **Privacy-Preserving ML** (PPML) or Privacy-Preserving Data Sharing and Analytics (PPDSA).

*“Privacy-preserving data sharing and analytics (PPDSA) methods **utilize cryptographic techniques**, which **inherently satisfy the confidentiality objective**.*

*The distinctive aspect of PPDSA approaches is their ability to achieve dissociability, **preventing authorized entities from establishing linkages between data and individuals' identities**, thereby enhancing privacy even with authorized data usage.*

Such technologies currently include, but are not limited to, secure multiparty computation, homomorphic encryption, zero-knowledge proofs, federated learning, secure enclaves, differential privacy, and synthetic data generation tools.”



NATIONAL STRATEGY TO ADVANCE PRIVACY-PRESERVING DATA SHARING AND ANALYTICS

A Report by the

FAST-TRACK ACTION COMMITTEE ON ADVANCING
PRIVACY-PRESERVING DATA SHARING AND ANALYTICS
NETWORKING AND INFORMATION TECHNOLOGY
RESEARCH AND DEVELOPMENT SUBCOMMITTEE

of the

NATIONAL SCIENCE AND TECHNOLOGY COUNCIL

March 2023

Privacy Preserving Machine Learning (PPML) @ SciPy 2023

Privacy guarantees are **the** most crucial requirement when it comes to analyse sensitive data. These requirements could be sometimes very stringent, so that it becomes a real barrier for the entire pipeline. Reasons for this are manifold, and involve the fact that data could not be shared, nor moved from their sites of resident. Let alone analyzed in

Northeastern University
College of Engineering

About Academics & Experiential Learning Research

NEWS & EVENTS / NEWS / MAKING AI MORE SECURE WITH PRIVACY-PRESERVING MACHINE LEARNING

Making AI More Secure with Privacy-Preserving Machine Learning

July 6, 2023

The Fourth AAAI Workshop on Privacy-Preserving Artificial Intelligence (PPAI-23)



February 13, 2023
© AAAI is an acronym event at Walter S. Washington Convention Center - Washington, DC, USA, Room 141A
PPAI will also be broadcasted at [aai.org](https://www.aai.org)

Event | June 2022

Apple Privacy-Preserving Machine Learning Workshop 2022



Earlier this year, Apple hosted the Workshop on Privacy-Preserving Machine Learning (PPML). This virtual event brought Apple and members of the academic research communities together to discuss the state of the art in the field of privacy-preserving machine learning through a series of talks and discussions over two days.

In this post we will introduce a new dataset for community benchmarking in PPML, and share highlights from workshop discussions and recordings of select workshop talks.

amazon | science

Research areas Blog News and features Publications Conferences

Privacy-preserving machine learning

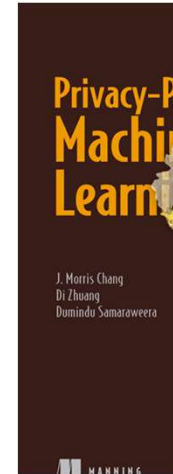
Microsoft | Research Our research Programs & events Blogs & podcasts About Sign up: Research Newsletter

Return to Blog Home
Microsoft Research Blog

Privacy Preserving Machine Learning: Maintaining confidentiality and preserving trust

Published November 9, 2021

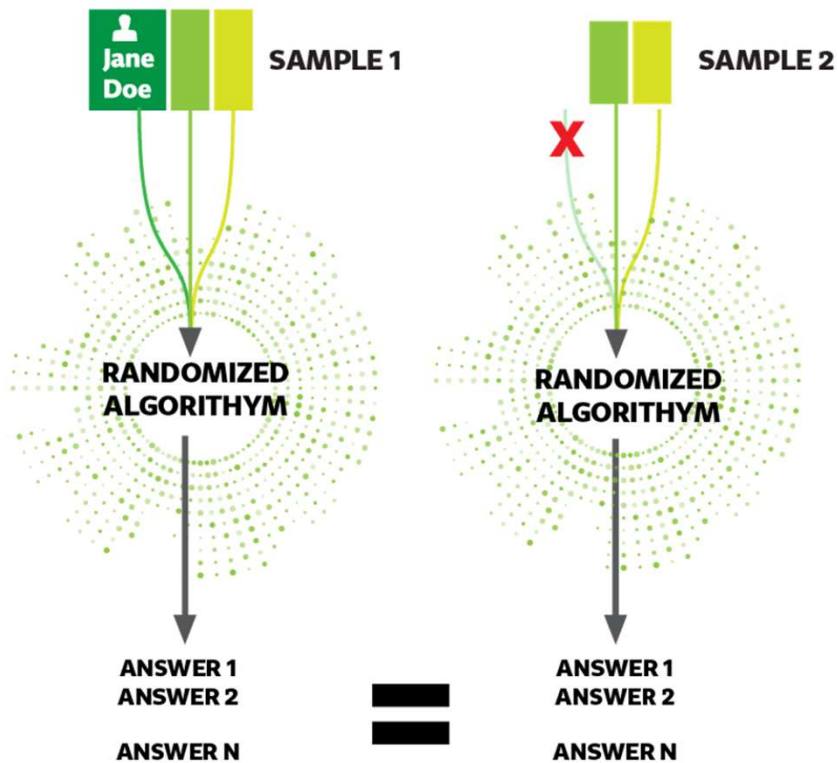
By [Victor Barthe](#), Principal Research Manager; [Robert Sim](#), Principal Research Manager; [Srinivas Aravamudan](#), Sr Principal Research Manager; [Noah Barth](#), Clarendon, Principal Researcher; [Malissa Chase](#), Principal Researcher; [David Jozef](#), Senior Applied Researcher; [Eduardo Lora](#), Principal Research Manager; [Boris Frenk](#), Principal Researcher; [Janice Struss](#), Chief Scientist & Technical Fellow; [Jan Delsing](#), Technical Fellow; [Sara Shalunov](#), Partner Director AI and Applied Research



The Second IACR School on Privacy-Preserving Machine Learning 2023

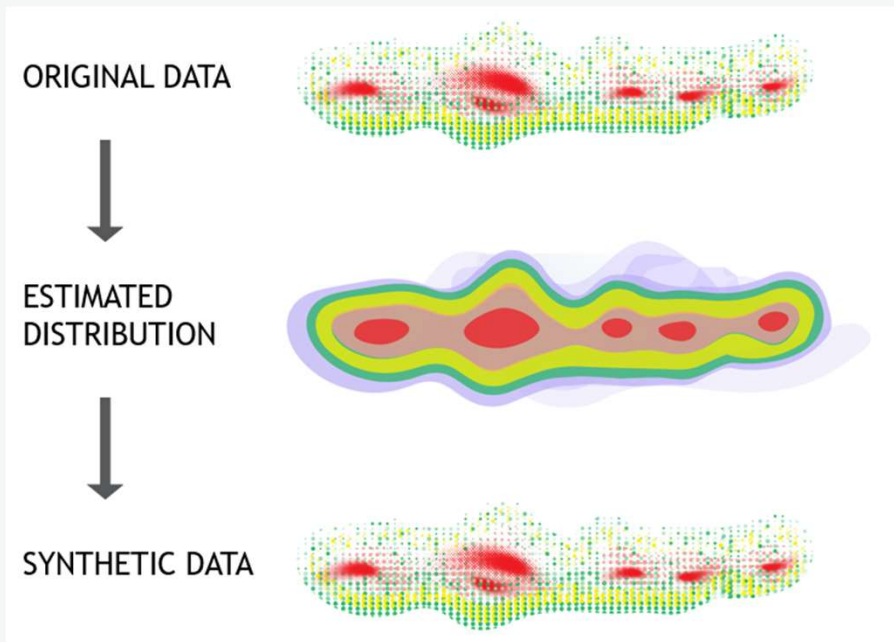
PRIVACY PRESERVING MACHINE LEARNING

DIFFERENTIAL PRIVACY



- Mathematical definition for privacy-preserving data analysis.
- **Goal:** Learn as much as possible about a group while learning as little as possible about any individual who is part of it.
- Outcome of analysis is essentially equally likely, independent of whether any individual joins, or refrains from joining, the dataset (Dwork, 2017).
- Achieved by adding “noise” (randomized responses) to data set.

SYNTHETIC DATA

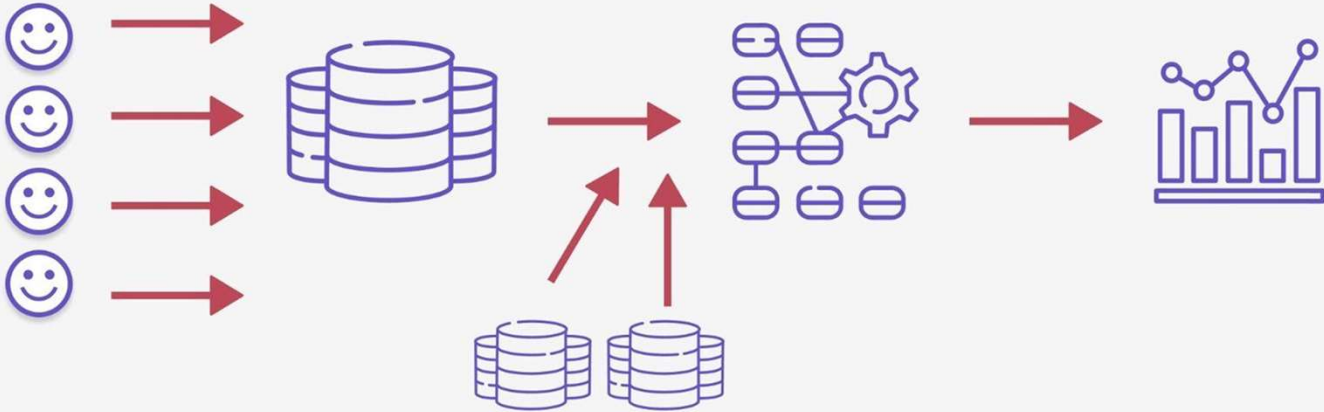


- Artificially generated data by an algorithm trained on a real data set.
- Goal: Preserving privacy in testing systems or training data for machine learning algorithms.
- Replaces original data while reproducing the statistical properties and patterns of the original set.

Collecting data privately
Secure aggregation
Local differential privacy

Computing on data privately
Homomorphic encryption
Confidential computing

Sharing data privately
Differential privacy



Joining data privately
Multi-party computation
Confidential computing

PERMANENT PROGRESS

Some examples:

Privacy rights (right to be deleted) & IP: Machine unlearning.

Hallucination: Retrieval Augmented Generation (RAG) - supplementing prompts with external data from an external data source (internet, APIs, databases, or documents).

Explainability/Interpretability: Decomposing groups of neurons into interpretable features... (*Anthropic*)

Model/Inference/Prompt Confidentiality: Trusted Execution Environments (TEEs) / Confidential computing

Best practices: Du-duplication, auditing, red teaming..

The collage features several elements:

- Top Left:** A snippet of a research paper titled "IN-CONTEXT UNLEARNING: LANGUAGE MODELS AS FEW SHOT UNLEARNERS" by Martin Pavlovic*, Seth Neel*, and Himanshu Lakkaraju*. The abstract discusses machine unlearning methods for LLMs.
- Top Right:** A diagram titled "EXISTING ENCRYPTION" comparing "Data at Rest" (encrypted in blob storage, database, etc.) with "Data in Transit" (encrypted between untrusted public or private networks). Below it, a red box labeled "NEW" shows "In use" data (Protect/Encrypt data that is in use, while in RAM and during computation.) with a source attribution to Confidential Computing Consortium (CCC).
- Middle Left:** A vertical label "579v2 [cs.LG] 12 Oct 2023".
- Bottom Section:** Three diagrams illustrating machine learning processes:
 - Rollout:** A flow from "Query" ("This movie is") to "LM" (GPT-2) to "Response" ("really great!").
 - Evaluation:** A flow from "Query + Response" ("This movie is really great!") to "Reward model" (Classifier/Rule/Human) to "Reward" (1.0).
 - Optimization:** A complex flow starting with "Query + Response" ("This movie is really great!") feeding into both an "Active model" and a "Reference model" (both LMs). The Active model outputs "log-probs", which are compared with the Reference model's "log-probs" to calculate "KL-div". This, along with "Reward", is used by "PPO" (Proximal Policy Optimization). A feedback loop labeled "Policy gradients optimize model" returns from PPO to the Active model.



OUTLOOK: RESPONSIBLE BY DESIGN

- Current AI systems will likely not meet all regulatory requirements.
- Attention will shift from merely understanding IF models comply with regulation to understanding how to **BUILD** models that comply with regulation.
- Harnessing Advanced ML with Transparency and Explainability



Thank you for your attention!

Follow me on LinkedIn and get in touch for learning more about Daiki - AI governance & enablement