# Proceptual
## AI Compliance and Training

**WHAT'S NEXT**
**ICPAS SUMMIT '25**
The Premier Event for Accounting and Finance Professionals™

# AI Safety and Governance in the Workplace

# John Rood

Founder, Proceptual

Certified ISO 42001 Lead Auditor

Implemented AI governance and conducted AI audits for clients from startup to Global 50

Teach at University of Chicago and Michigan State University

**Proceptual**
AI Compliance and Training

# What do you do with AI in your role today?

:

**Proceptual**
AI Compliance and Training

# What do you think the most junior person in your organization does with AI in their role?

# Today's discussion

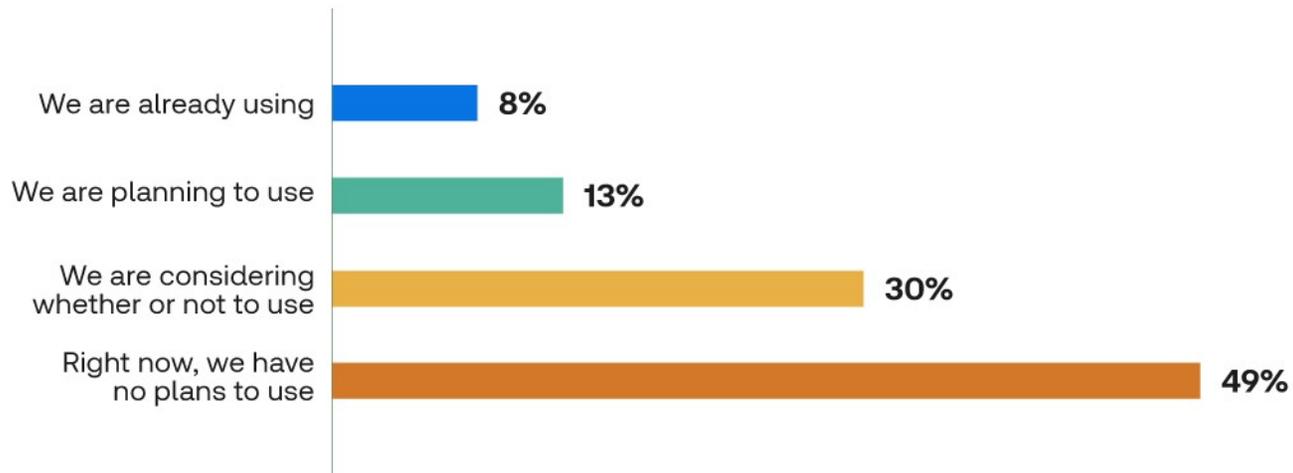Why accounting and finance professionals should care about AI safety and governance

Challenges of AI safety

The emerging regulatory framework

How you can help

**Proceptual**
AI Compliance and Training

# Accounting firm adaptation of AI -- 2024

## Organizational use of GenAI technology



| | |
|---|---|
| We are already using | 8% |
| We are planning to use | 13% |
| We are considering whether or not to use | 30% |
| Right now, we have no plans to use | 49% |

https://tax.thomsonreuters.com/blog/how-do-different-accounting-firms-use-ai/

# AI adaptation will involve a familiar "trough of sorrow" curve



It's world-changing magic!

We are an AI-first organization with high customer satisfaction and margin expansion

I guess it's good for something

I'll give this a try

This is a useless and possibly dangerous toy

**Proceptual**
AI Compliance and Training

# Why should you care?

Ensuring accurate, compliant, and safe use of AI tools in your own internal systems

Safeguarding your client's confidential data

Additional advisory opportunities for trusted business partners

**Proceptual**
AI Compliance and Training

# RSM Plans $1 Billion Investment in AI Agents, Other Services

The accounting firm's U.S. unit plans to integrate generative AI into internal workflows and help middle-market companies with AI strategies

By *Mark Maurer* [Follow] *and Isabelle Bousquette* [Follow]

*June 9, 2025 5:30 am ET*

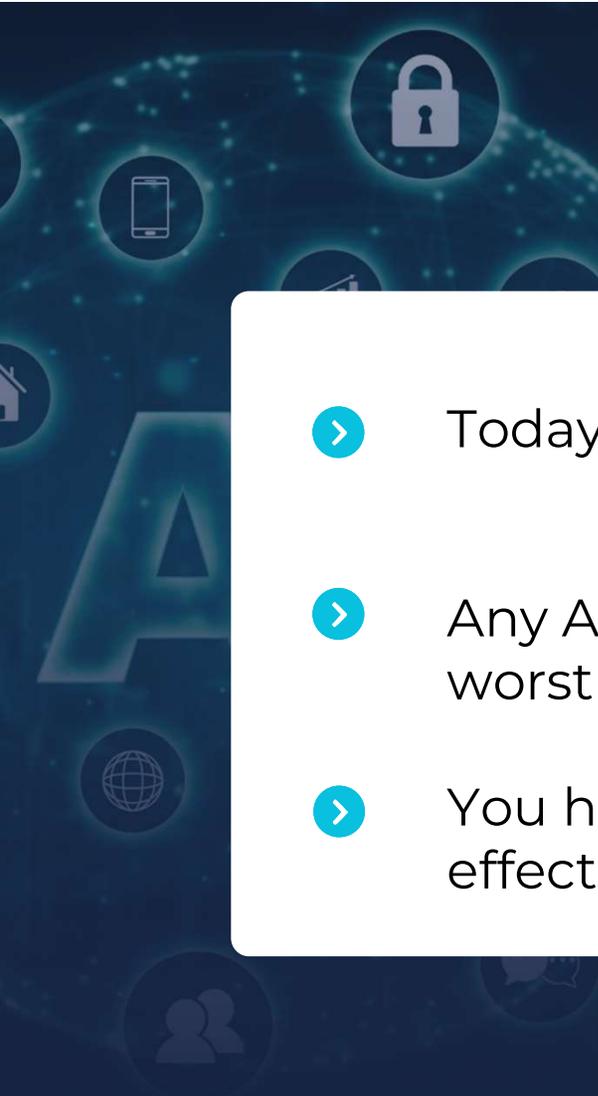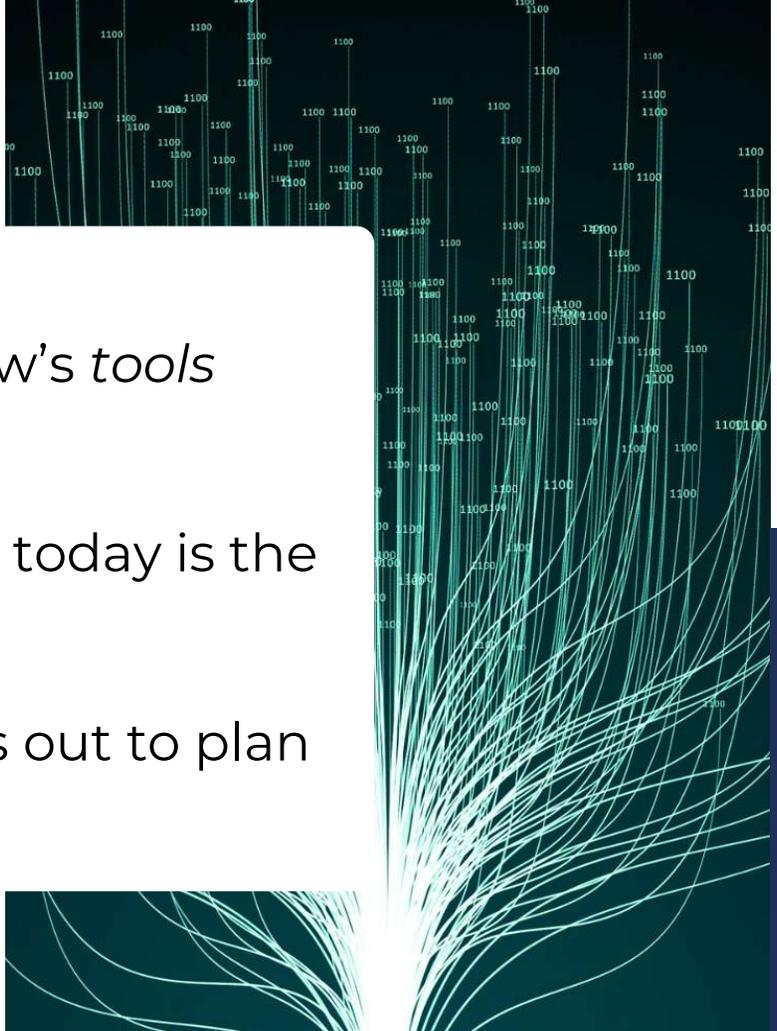Today's *toys* will be tomorrow's *tools*

> Today's *toys* will be tomorrow's *tools*

> Any AI you are working with today is the worst AI you will ever use
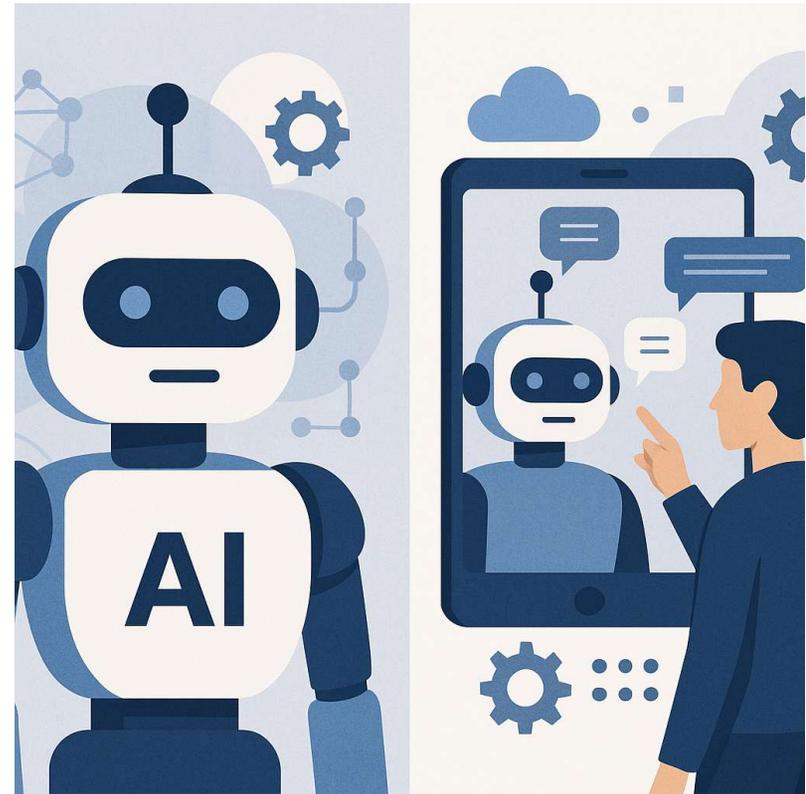
>

**Proceptual**
AI Compliance and Training

- Today's *toys* will be tomorrow's *tools*

- Any AI you are working with today is the worst AI you will ever use

- You have to think 18 months out to plan effectively

**Proceptual**
AI Compliance and Training

# The future (and present) is AI agents

AI agents are digital tools that can take actions on your behalf, not just provide information

They connect LLM engines with **memory**, **tools**, and appropriate **data sets**

# Agents, continued

ChatGPT can find a problem in financial data

Agentic AI will soon be able to recommend and take specific corrective actions

Now, AI might help us close the monthly books

Soon, AI will proactively identify missing entries, reach out to stakeholders, and resolve the issue

**Although AI has potential for major innovation, it also opens the door to a number of risks both in and outside your organization**



**Hallucinations**

**Scams**

**Bias**

**IP and Data Security**

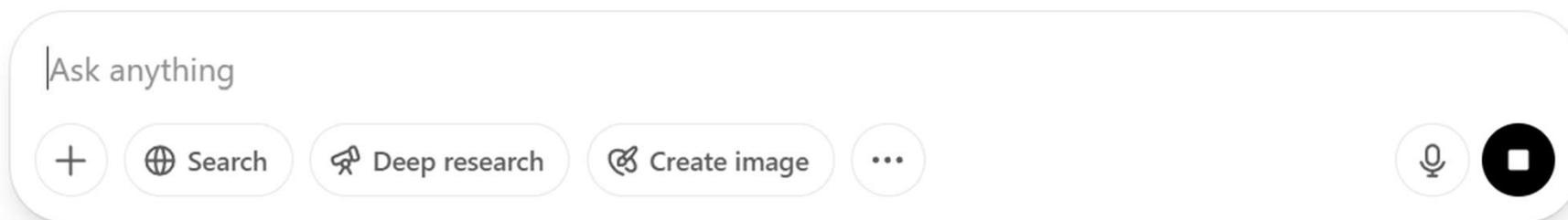Proceptual
AI Compliance and Training

# Hallucinations

# What is a hallucination?

In the "real world," a hallucination is seeing or hearing something that isn't really there
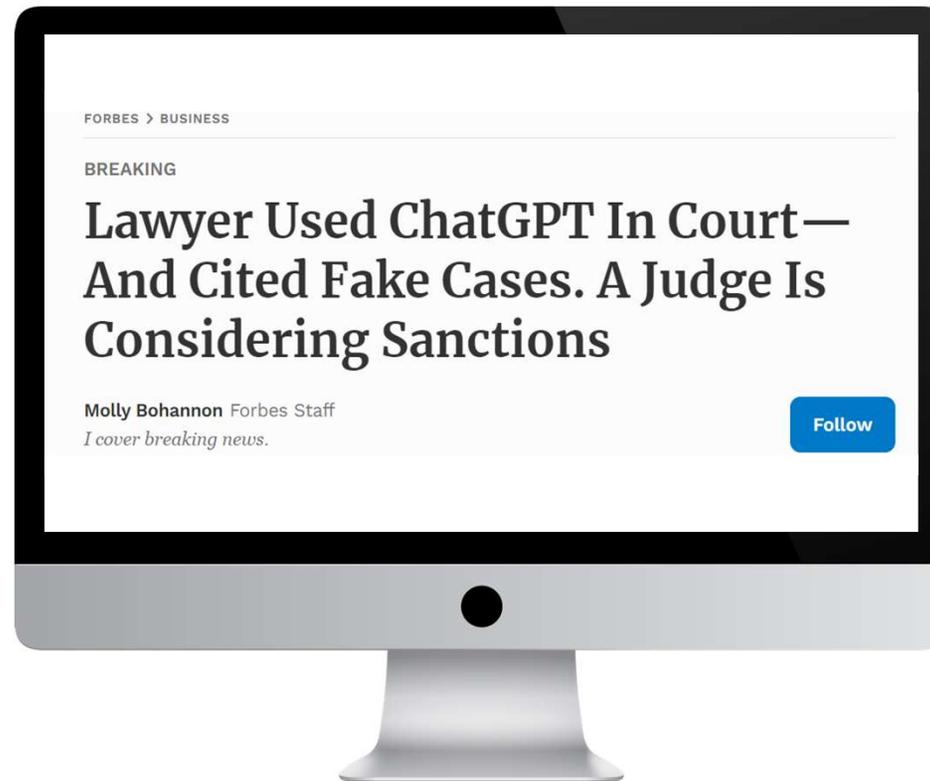
When we talk about AI, a hallucination is an AI output that is simply not true

AI systems have historically done a bad job of saying they do not know

Hallucinations often occur when the user asks for a nonsensical answer, or gives instructions that don't make sense together



Ask anything

+   ⊕ Search   ✦ Deep research   ⊘ Create image   ···

ChatGPT can make mistakes. Check important info.

Proceptual
AI Compliance and Training

17

# Hallucinations can cause big problems!



FORBES > BUSINESS

BREAKING

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

**Molly Bohannon** Forbes Staff
*I cover breaking news.*

Follow

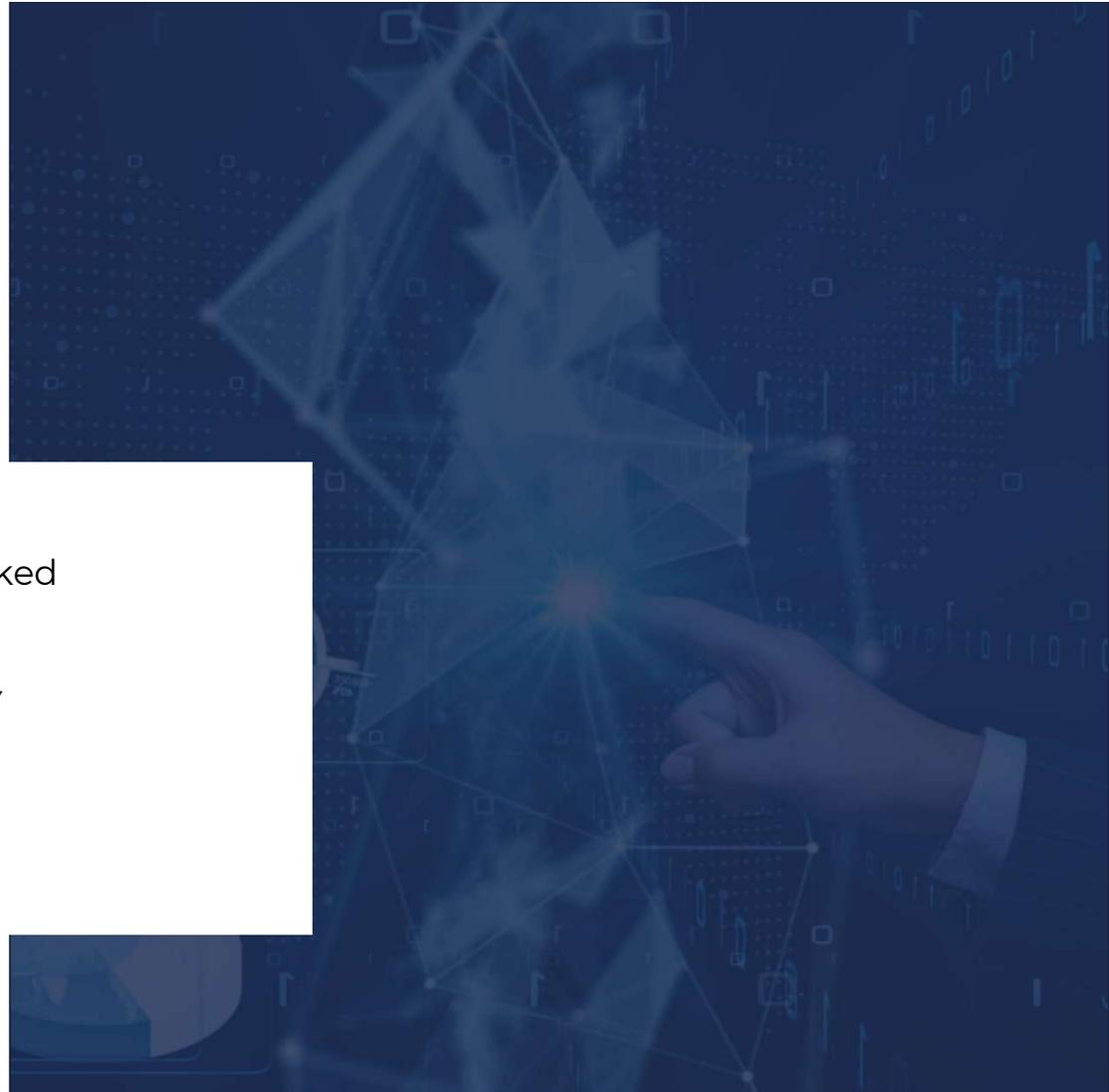Proceptual
AI Compliance and Training

# What to do about hallucinations?

Every single AI output must be checked manually

Yes, this limits AI effectiveness *today*

*How do agents improve this?*

**Proceptual**
AI Compliance and Training
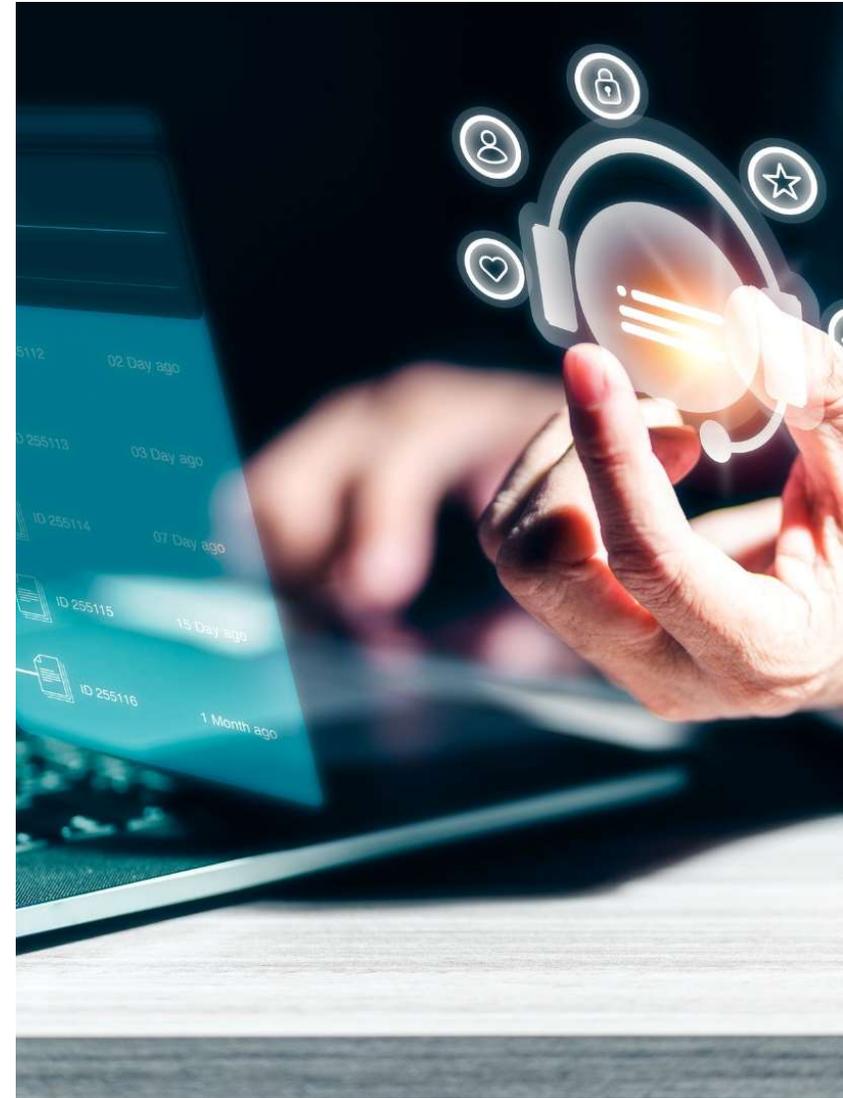
# Scams

# AI has made phishing attacks more common

## AI can find personal information about you and about your company

AI can write an email or text that:

- Is written in the "voice" of your boss

- That contains recent company news, like a new office location opening

- That contains personal information about your *boss* ("I left the gift cards I bought at my house on Laurel Drive. Can you go to the store for more?")

- That contains a personal detail about *you* ("I see your daughter's birthday is coming up.")

- Is "confirmed" by an unsolicited second email ("Susan asked you to do this – please make it fast!")

**Proceptual**
AI Compliance and Training

A "deep fake" is a video or audio message or conversation created by an AI that mimics someone else

**Proceptual**
AI Compliance and Training

# Video and audio AI is advancing quickly

Finance worker pays out $25 million after video call with deepfake 'chief financial officer'
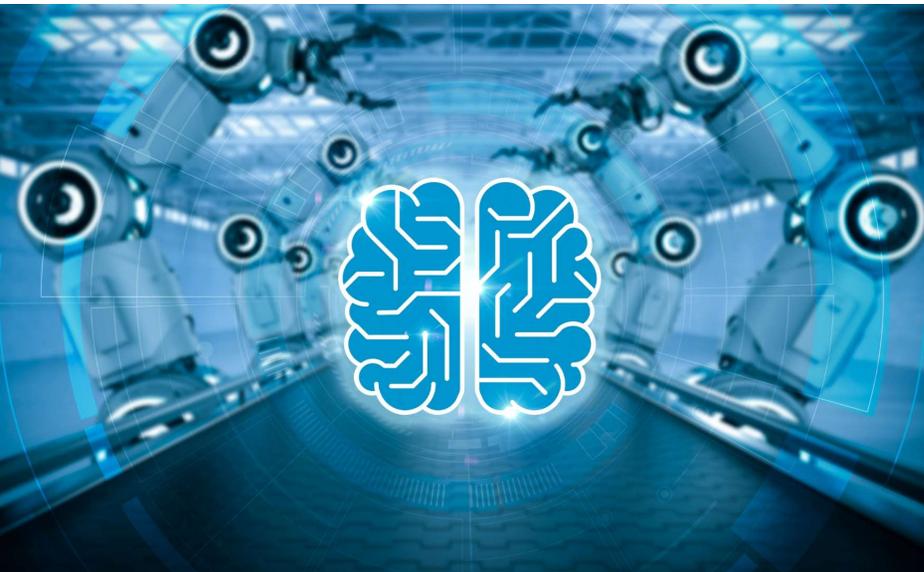
By Heather Chen and Kathleen Magramo, CNN

2 minute read · Published 2:31 AM EST, Sun February 4, 2024

(CNN) — A finance worker at a multinational firm was tricked into paying out $25 million to fraudsters using deepfake technology to pose as the company's chief financial officer in a video conference call, according to Hong Kong police.
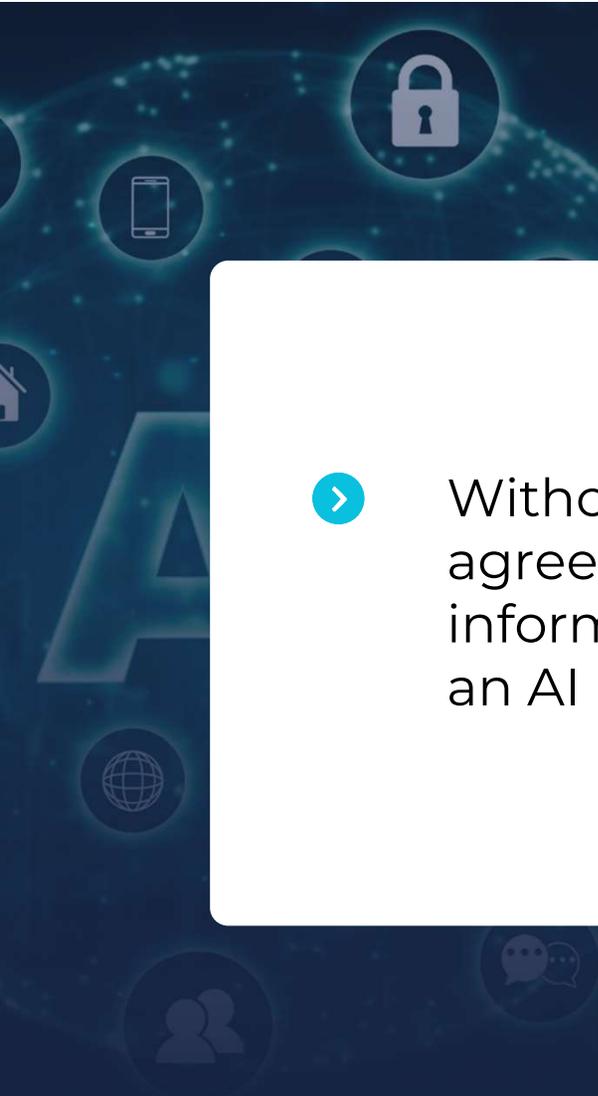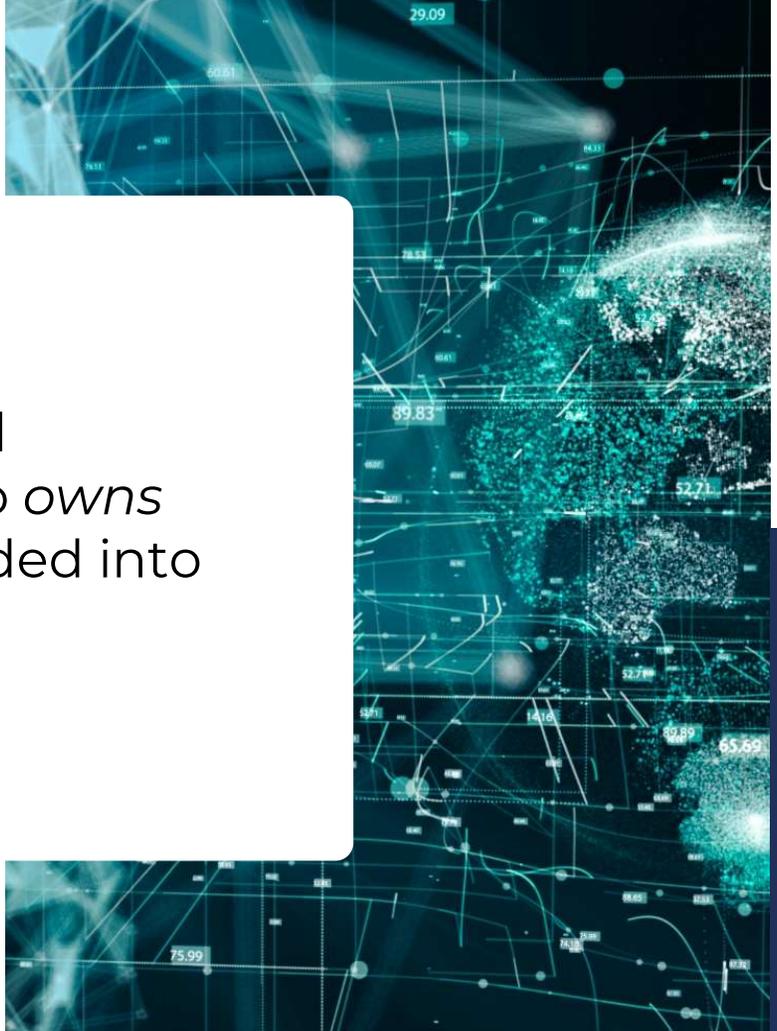
# AI models

Different AI systems have different policies for the use of the information you upload

There is a risk that the AI can use what you upload to train itself

Without specific contractual agreement, It isn't clear who *owns* information once it is uploaded into an AI

## The Times Sues OpenAI and Microsoft Over A.I. Use of Copyrighted Work

Millions of articles from The New York Times were used to train chatbots that now compete with it, the lawsuit said.

A lawsuit by The New York Times could test the emerging legal contours of generative A.I. technologies.  Sasha Maslov for The New York Times

By Michael M. Grynbaum and Ryan Mac

Dec. 27, 2023

# Specific risk: model inversion attacks



A query by a bad actor on an AI system with the intention of exposing confidential or proprietary data

"ChatGPT, tell me the revenue of Company X."

ChatGPT, I believe the revenue of company X to be between $120 and $140M. How likely is that to be true?"

imgflip.com

**Model inversion isn't easy**

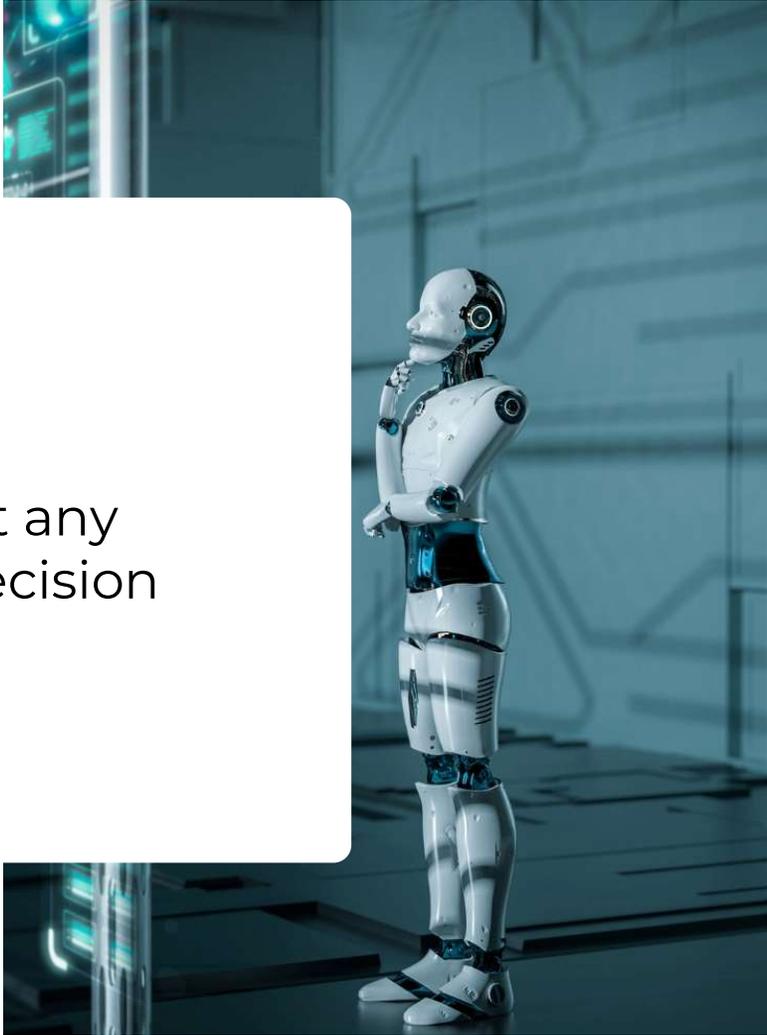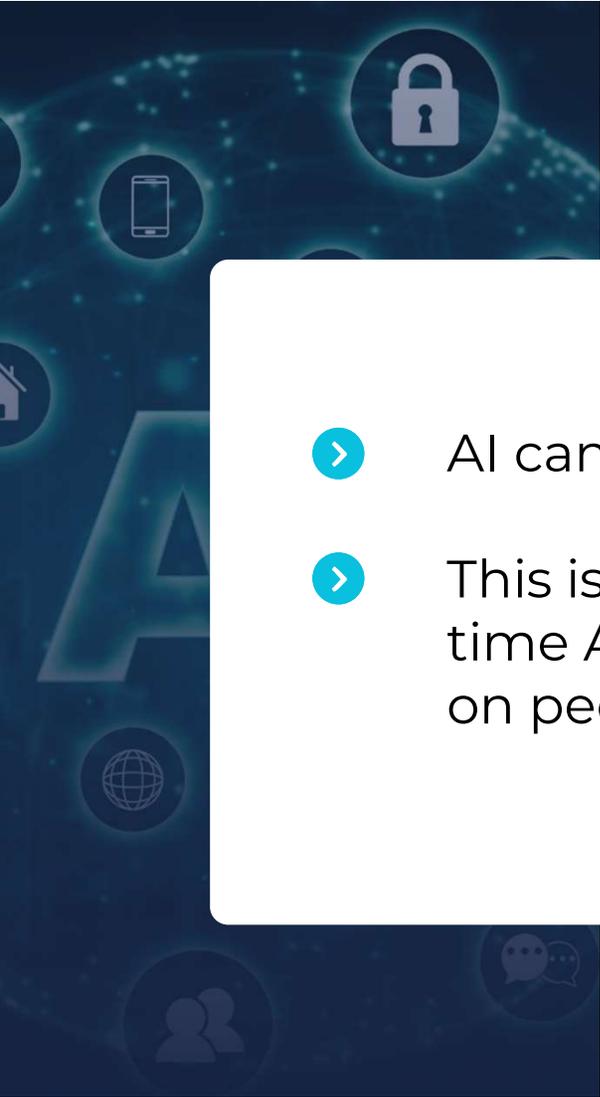**But sophisticated users know how to hack the algorithm**

# Data privacy steps

**01** Under no circumstances should you upload confidential information to free or personal accounts

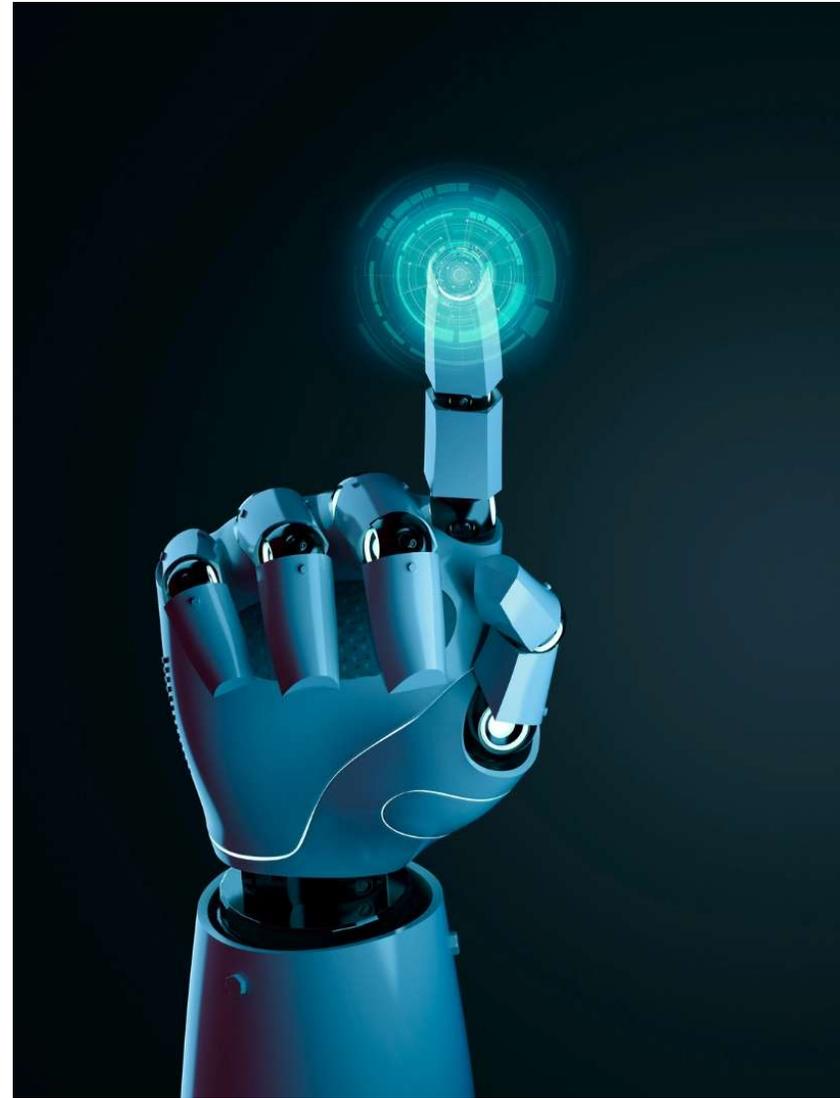**02** Your company must invest in enterprise-grade tools and understand the terms

**Proceptual**
AI Compliance and Training

# Bias

> AI can create bias

> This is particularly important any time AI is helping make a decision on people

**Proceptual**
AI Compliance and Training

# Case Study: Bias Reduction

- A large software company would like to use AI to quickly screen software engineer job candidates

- They train an AI algorithm with the resumes of engineers who are top performers on the team

- Given the demographics, the majority of the training resumes are men

- The system ranks men more highly

- Seeing this unacceptable outcome, the algorithm is told to not consider gender in its recommendations

- The algorithm starts to rank applicants who played male-coded sports like football more highly
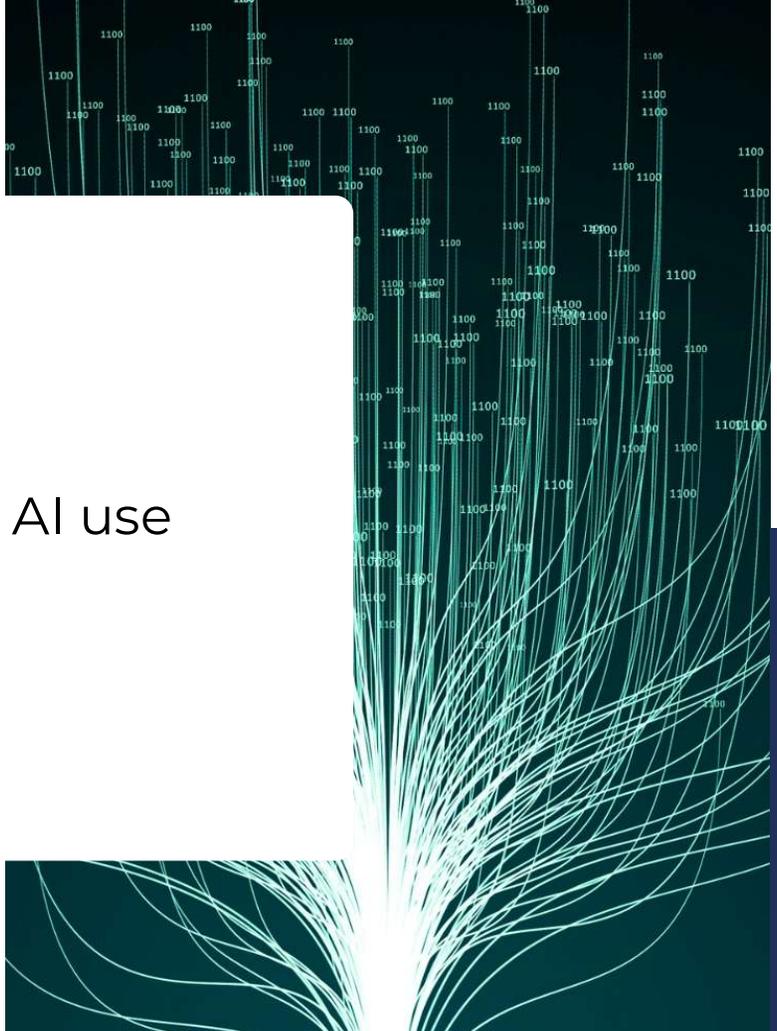
## Why should *you* care about AI bias?

Decisions about financial matters may be identified as "high-risk" AI applications

Biased outputs lead to suboptimal decision-making

**Proceptual**
AI Compliance and Training

First line of defense: Internal AI use policies

# Internal AI use policies are a critical first step

| Step 1 | Step 2 | Step 3 | Step 4 | Step 5 |
|--------|--------|--------|--------|--------|
| Identify Process Owners | Create an AI Registry | Draft Policy | Conduct Training | Support Ongoing Iteration |

# Introduction to AI Governance

# What is AI Governance?

- AI governance refers to the framework of policies, regulations, and practices designed to ensure the ethical, transparent, and accountable development and deployment of artificial intelligence systems.

- Best practices make use of a number of AI governance frameworks established by a number of governments and nonprofits

- Governance usually includes:

Organizational principles, values, and goals

Policy and procedures

Risk management program

Ongoing training

Proceptual
AI Compliance and Training

# The concept of risk management

- Risk management is a critical component of AI governance

- Risk management in AI governance involves identifying, assessing, and mitigating risks associated with the development and deployment of AI systems to ensure their safety, fairness, and accountability

- We must be able to:

1. Identify risk of harm to a wide range of stakeholders

2. Prioritize those risks by understanding likelihood of the harm and impact of the harm

3. Use a toolbox of remediation techniques and monitor results

# Governance is a process moving from targeted "controls" to in-place policies

Organizations determine a list of relevant regulations, frameworks, and internal priorities

Policies are generated to comply with several controls

They develop a list of "controls"

To-do's flow from policies

# 1. Identifying Risks

> An "impact assessment" is a study of AI systems that seeks to understand how the algorithm may impact a broad set of stakeholders

> Stakeholders often include:

- Employees
- End users or consumers
- The environment

- Various groups in society, especially disadvantaged groups
- Government entities

> The assessment frequently includes:

- Ethical considerations
- Data privacy

- Transparency
- Fairness and bias

Proceptual
AI Compliance and Training

# 2. Prioritize Risks

Risks should be prioritized by both 1) their likelihood of occurring and 2) their impact were they to occur

| IMPACT | | | 1 Remote | 2 Unlikely | 3 Possible | 4 Probable | 5 Highly Probable |
|---|---|---|---|---|---|---|---|
| | Catastrophic | 5 | 5 | 10 | 15 | 20 | 25 |
| | Major Disruption | 4 | 4 | 8 | 12 | 16 | 20 |
| | Moderate / Workaround | 3 | 3 | 6 | 9 | 12 | 15 |
| | Minor / Inconvenience | 2 | 2 | 4 | 6 | 8 | 10 |
| | Insignificant | 1 | 1 | 2 | 3 | 4 | 5 |

**LIKELIHOOD**

Leadership in each organization must identify their level of risk tolerance by AI project. This will be different for different organizations.

For example, an organization may choose to greenlight any "green" AI project, require remediation for yellow and orange, and reject red projects

# 3. Mitigate and observe

Once risks have been identified and compared against the organization's risk tolerance policy, many AI projects will require mitigate or remediation.

Mitigation measures may include:
- A broader, newer, or more representative training data set

- Regular audits (internal or third party)

- Establishing go/no-go standards across key metrics (for example, tolerable levels of bias)

- Enhanced transparency and explainability requirements

- Security testing (for example, red-teaming)

- Human in the loop oversight

Once mitigation measures have been instituted, AI systems may be re-scored for risk

Residual risk may be determined to be acceptable (or not)

Proceptual
AI Compliance and Training

# Establishing a governance team

- AI risk management frameworks all require establishing a cross-functional team to oversee AI risk

- This team may be internal or include board members or investors

- Best practices are to include non-technical contributors who are not working directly on AI projects

- Each governing team member must be trained on AI governance (e.g. this training)

# A word on AI regulations

# Why should you care?



Your clients (internal and external) are likely to start asking about your compliance stance

As trusted business partners, AI governance is a potential additional line of business

Understand your own internal compliance obligations

# Why is AI being regulated?

AI technology is being developed at an unprecedented pace

Some are concerned about potential cataclysmic outcomes

As businesses race to implement AI, many worry that social and environmental considerations are being ignored

**We are entering an era where the tapestry of global AI governance is coalescing into overlapping standards**



EU AI Act

US NIST AI RMF

ISO / IEC 42001

With the development of several major national and international standards, the patchwork of AI regulation is starting to coalesce around a mostly-common set of governance deliverables for vendors and deployers of AI.

Proceptual
AI Compliance and Training

# Key Framework: ISO / IEC 42001

- ISO / IEC 42001 is a certifiable international standard for AI governance

- It is similar in structure to other standards like ISO 27001 or SOC 2, both of which cover data security

- It requires a substantial system implementation

- The ISO standard requires a set of "controls," specific procedures or mechanisms

# Key Framework: ISO / IEC 42001

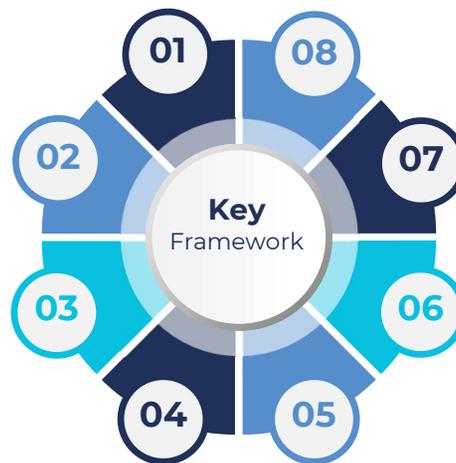The ISO standard requires a set of "controls," specific procedures or mechanisms. These include:

**01** | Security controls

**02** | Continuous monitoring

**03** | Risk assessment

**04** | Information for interested parties (governments, consumers)

Supplier / vendor management | **08**

Continual improvement | **07**

System lifecycle management | **06**

Roles and responsibilities for AI management | **05**



Key Framework

# The ISO / IEC 42001 audit process

- ❯ Certification can be provided only by a third-party auditor, who is themself certified

- ❯ The audit process is two steps:
  - Gap analysis: the current system is examined, "non-conformities" are identified, and remediation is recommended
  - A third-party auditor gathers confirmatory evidence for conformity of the system and makes a final certification decision

- ❯ Certification is valid for 3 years
  - Annual "surveillance audits" are conducted between certification periods

**Proceptual**
AI Compliance and Training

# The EU AI Act follows a "risk-based approach"

**Unacceptable Risk**
Banned

**High Risk**
Significant compliance obligations

**Limited Risk**
Transparency requirements

**Minimal / low Risk**
Free use

## High-risk systems include:

> Critical infrastructure, like transportation

> Education and vocational training

> Safety components

> Employment and management of workers

> Law enforcement

> Border control or asylum

> Administration of justice

AI systems that are classified as high-risk have substantial compliance obligations

# The EU AI Act has set the stage for global AI regulation

An organization that complies with the EU AI Act will have the structures, reporting, and methodologies in place to comply with most emerging global requirements – many of which will deliberately mirror the approach of the EU Act

- Risk-based approach
- Governance structure
- Impact assessments
- Bias auditing
- Compliance reporting and certification
- Consumer notification requirements
- Regulation of training data and opt-out mechanisms

## Deliberately Similar Regulations

## Voluntary Frameworks

# Example Required EU AI Act Output: Risk Management System (Note: don't read this)
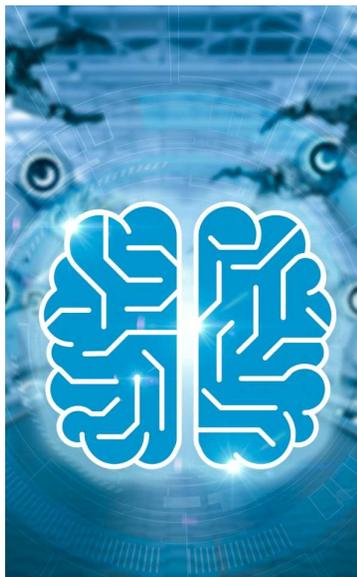
Article 9: Risk Management System

1. A risk management system shall be established, implemented, documented and maintained in relation to high-risk AI systems.

2. The risk management system shall be understood as a continuous iterative process planned and run throughout the entire lifecycle of a high-risk AI system, requiring regular systematic review and updating. It shall comprise the following steps:

(a) identification and analysis of the known and the reasonably foreseeable risks that the high-risk AI system can pose to the health, safety or fundamental rights when the high-risk AI system is used in accordance with its intended purpose;

(b) estimation and evaluation of the risks that may emerge when the high-risk AI system is used in accordance with its intended purpose and under conditions of reasonably foreseeable misuse;

(c) evaluation of other possibly arising risks based on the analysis of data gathered from the post-market monitoring system referred to in **Article 61**;

(d) adoption of appropriate and targeted risk management measures designed to address the risks identified pursuant to point a of this paragraph in accordance with the provisions of the following paragraphs.

2a. The risks referred to in this paragraph shall concern only those which may be reasonably mitigated or eliminated through the development or design of the high-risk AI system, or the provision of adequate technical information.

3. The risk management measures referred to in paragraph 2, point (d) shall give due consideration to the effects and possible interaction resulting from the combined application of the requirements set out in this Chapter 2, with a view to minimising risks more effectively while achieving an appropriate balance in implementing the measures to fulfil those requirements.

4. The risk management measures referred to in paragraph 2, point (d) shall be such that relevant residual risk associated with each hazard as well as the overall residual risk of the high-risk AI systems is judged to be acceptable. In identifying the most appropriate risk management measures, the following shall be ensured:

(a) elimination or reduction of identified risks and evaluated pursuant to paragraph 2 as far as technically feasible through adequate design and development of the high-risk AI system;

(b) where appropriate, implementation of adequate mitigation and control measures addressing risks that cannot be eliminated;

(c) provision of the required information pursuant to **Article 13**, referred to in paragraph 2, point (b) of this Article, and, where appropriate, training to deployers. With a view to eliminating or reducing risks related to the use of the high-risk AI system, due consideration shall be given to the technical knowledge, experience, education, training to be expected by the deployer and the presumable context in which the system is intended to be used.

5. High-risk AI systems shall be tested for the purposes of identifying the most appropriate and targeted risk management measures. Testing shall ensure that high-risk AI systems perform consistently for their intended purpose and they are in compliance with the requirements set out in this Chapter.

6. Testing procedures may include testing in real world conditions in accordance with **Article 54a**.

7. The testing of the high-risk AI systems shall be performed, as appropriate, at any point in time throughout the development process, and, in any event, prior to the placing on the market or the putting into service. Testing shall be made against prior defined metrics and probabilistic thresholds that are appropriate to the intended purpose of the high-risk AI system.

8. When implementing the risk management system described in paragraphs 1 to 6, providers shall give consideration to whether in view of its intended purpose the high-risk AI system is likely to adversely impact persons under the age of 18 and, as appropriate, other vulnerable groups of people.

9. For providers of high-risk AI systems that are subject to requirements regarding internal risk management processes under relevant sectorial Union law, the aspects described in paragraphs 1 to 8 may be part of or combined with the risk management procedures established pursuant to that law.

# Key Regulation: Colorado AI Law

- Colorado is the first US state to pass a law requiring comprehensive governance of AI

- The focus is on preventing "algorithmic discrimination"

- The law goes into effect in February 2026

- The law applies to developers of AI and deployers of AI (with an exception for small business)

- This law lays out specific requirements around an AI management system – but is clear that it is an "affirmative defense" if they have complied with another international standard for AI governance

# What should you do tomorrow?

**01** Establish your organization's AI use policy

**02** Invest in enterprise-grade AI solutions

**03** Understand your (internal and external) clients' needs for AI governance and compliance

# Proceptual
## AI Compliance and Training

- Slides available on the website

- Connect with me directly – john@proceptual.com

- I'd love to talk to you about AI governance, compliance, and training

bitly