



The *Next* 6,173 AI Incident Reports

Sean McGregor, Ph.D., ExecDir RAIC
Fellow of the Berkman Klein Center for
Internet & Society at Harvard University

NIST Workshop on AI Incident Management
May 14th, 2026

Our *First* 6,173 AI Incident Reports

Sean McGregor, Ph.D., ExecDir RAIC
Fellow of the Berkman Klein Center for
Internet & Society at Harvard University

NIST Workshop on AI Incident Management
May 14th, 2026

Let's look back to 2018...



Aviation Safety: Fatalities per trillion RPK



NTSB

National Transportation Safety Board

Investigations | Safety Research | News & Events | Advocacy | Family Assistance | About Us

Aviation Accident Database & Synopses

For cases after 2008, use [CAROL Query](#).
Learn about changes to our [search options](#).

The NTSB aviation accident database contains information from 1962 and later about civil aviation accidents and selected incidents within the United States, its territories and possessions, and in international waters. Generally, a preliminary report is available online within a few days of an accident. Factual information is added when available, and when the investigation is completed, the preliminary report is replaced with a final description of the accident and its probable cause. Full narrative descriptions may not be available for dates before 1993, cases under revision, or where NTSB did not have primary investigative responsibility.

- [Monthly lists](#) - accidents sorted by date, updated daily.
- [Downloadable datasets](#) - one complete dataset for each year beginning from 1982, updated monthly in Microsoft Access 2000 MDB format; this site also provides weekly "change" updates and complete documentation.
- [GILS record](#) - complete description of the accident database, including definition of "accident" and "incident".
- [FAA incident database](#) - complete information about incidents, including those not investigated by NTSB, is provided by the Federal Aviation Administration.

[Help](#)

Accident/Incident Information

Event Start Date (mm/dd/yyyy)

Aircraft

Category

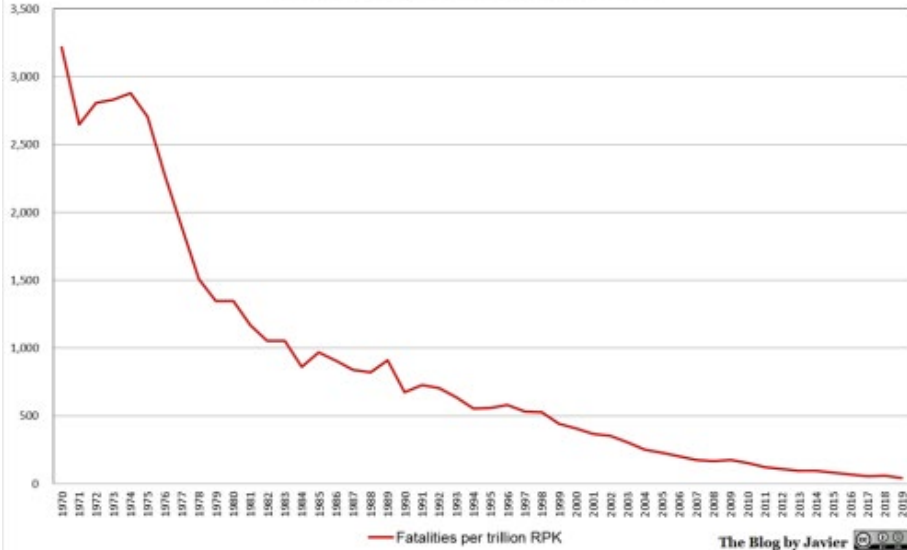
All

Amateur Built

All



Aviation Safety: Fatalities per trillion RPK



The screenshot shows the NTSB website's interface for the Aviation Accident Database. It includes a navigation menu with links for Investigations, Safety Research, News & Events, Advocacy, Family Assistance, and About Us. The main heading is "Aviation Accident Database & Synopses". Below this, there is a note for cases after 2008 to use CAROL Query and a link to search options. A paragraph describes the database's scope and the availability of preliminary and final reports. A list of resources is provided, including monthly lists, downloadable datasets, GILS records, and the FAA incident database. There is also a "Help" link and a search form with fields for "Event Start Date", "Event End Date", "Category", and "Amateur Built".

This is what we want for AI

Let's look back to 2018...

Let's look back to 2018...

The image shows a screenshot of a Microsoft Excel spreadsheet. The spreadsheet is filled with a grid of data, including text, numbers, and dates. The columns are labeled with dates from 2018, and the rows contain various entries. A blue selection box is visible over a portion of the grid, highlighting a specific area of data. The spreadsheet appears to be a detailed record or log, possibly related to a project or event in 2018.

A humble spreadsheet

Let's look back to 2018...

Year	Month	Day	Time	Location	Event	Description	Remarks
2018	Jan	1	08:00
2018	Jan	2	08:00
2018	Jan	3	08:00
2018	Jan	4	08:00
2018	Jan	5	08:00
2018	Jan	6	08:00
2018	Jan	7	08:00
2018	Jan	8	08:00
2018	Jan	9	08:00
2018	Jan	10	08:00
2018	Jan	11	08:00
2018	Jan	12	08:00
2018	Jan	13	08:00
2018	Jan	14	08:00
2018	Jan	15	08:00
2018	Jan	16	08:00
2018	Jan	17	08:00
2018	Jan	18	08:00
2018	Jan	19	08:00
2018	Jan	20	08:00
2018	Jan	21	08:00
2018	Jan	22	08:00
2018	Jan	23	08:00
2018	Jan	24	08:00
2018	Jan	25	08:00
2018	Jan	26	08:00
2018	Jan	27	08:00
2018	Jan	28	08:00
2018	Jan	29	08:00
2018	Jan	30	08:00
2018	Jan	31	08:00

What do we call these?

A humble spreadsheet

AID

What do we call these?

accidents, failures, malfunctions?

NID

What do we call these?

accidents, failures, malfunctions?

 **Doesn't cover intended harms**

NID

What do we call these?

accidents, failures, malfunctions?



Doesn't cover intended harms

vulnerabilities, threats, hazards

NID

What do we call these?

accidents, failures, malfunctions?



Doesn't cover intended harms

vulnerabilities, threats, hazards



Not event-y

AID

What do we call these?

accidents, failures, malfunctions?



Doesn't cover intended harms

vulnerabilities, threats, hazards



Not event-y

catastrophes, misuse, recalls, crashes

AID

What do we call these?

accidents, failures, malfunctions?

👎 **Doesn't cover intended harms**

vulnerabilities, threats, hazards

👎 **Not event-y**

catastrophes, misuse, recalls, crashes

👎 **Too specific**

AID

What do we call these?

accidents, failures, malfunctions?

👎 **Doesn't cover intended harms**

vulnerabilities, threats, hazards

👎 **Not event-y**

catastrophes, misuse, recalls, crashes

👎 **Too specific**

Incidents?

AID

What do we call these?

accidents, failures, malfunctions?

👎 **Doesn't cover intended harms**

vulnerabilities, threats, hazards

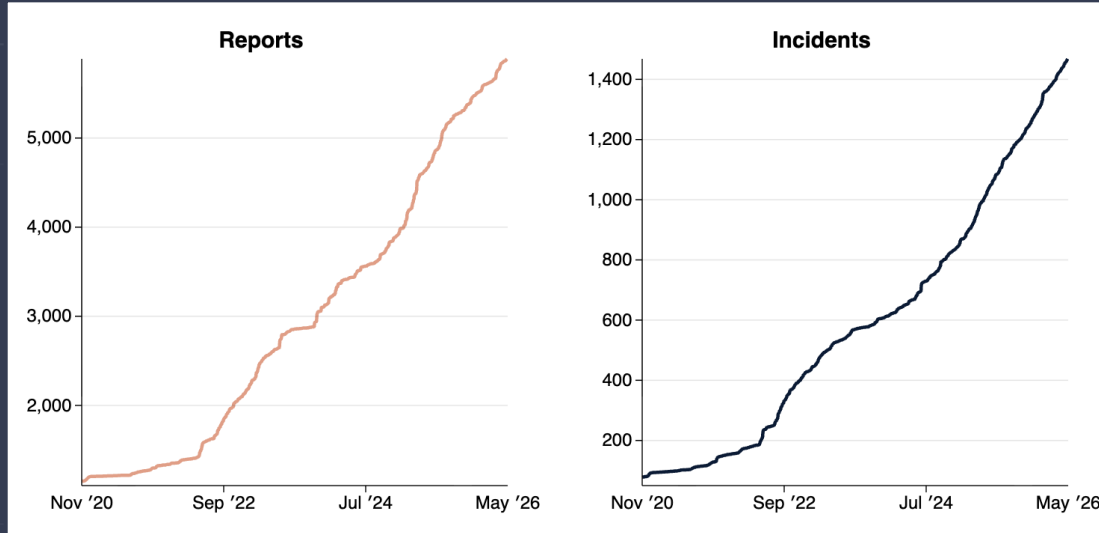
👎 **Not event-y**

catastrophes, misuse, recalls, crashes

👎 **Too specific**

Incidents?

Yes! Incidents ⚡



6,173 AI Incident Reports
(and counting)

Discover

Submit

Submitted Incident Report List

The following incident reports have been **submitted** by users and are pending review by editors. Only editors may promote these records to incident reports in the database.

TITLE	SUBMITTERS	DATES	EDITORS	STATUS	ACTIONS
<input type="checkbox"/> Search 10 records...	All ▾	Incident date ▾	All ▾	All ▾	
<input type="checkbox"/> Her daughter was unraveling, and she didn't know why. Then she found the AI chat logs.	Ashmita Rajmohan	Sub: 2026-01-10 Pub: 2025-12-23	Daniel Atherton	In Review	<input type="button" value="Review"/> <input type="button" value="Claim"/>
<input type="checkbox"/> LLMs have a misalignment feature that uses strategic deception	Anonymous	Inc: 2026-04-25 Sub: 2026-04-25 Pub: 2026-04-25	Daniel Atherton	In Review	<input type="button" value="Review"/> <input type="button" value="Claim"/>
<input type="checkbox"/> Claude-powered AI coding agent deletes entire company database in 9 seconds – backups zapped, after Cursor tool powered by Anthropic's Claude goes rogue	Anonymous	Inc: 2026-04-25 Sub: 2026-04-27 Pub: 2026-04-27	Daniel Atherton	In Review	<input type="button" value="Review"/> <input type="button" value="Claim"/>
<input type="checkbox"/> AI Coding Agent Powered by Claude Opus 4.6 Deletes Production Database in 9 Seconds	Atluxity	Sub: 2026-04-29 Pub: 2026-04-28	Daniel Atherton	In Review	<input type="button" value="Review"/> <input type="button" value="Claim"/>
<input type="checkbox"/> Google Gemini drafted immigration email that caused loss of university	Anonymous	Inc: 2026-03-03 Sub: 2026-05-01	Daniel Atherton	In Review	<input type="button" value="Review"/> <input type="button" value="Claim"/>

Welcome to the AIID

Discover Incidents

Spatial View

Table View

List view

Entities

Taxonomies

Submit Incident Reports

Submission Leaderboard

Blog

AI News Digest

Risk Checklists

Random Incident

Account



Discover

Submit

Submitted Incident Report List

The following incident reports have been **submitted** by users and are pending review by editors. Only editors may promote these records to incident reports in the database.

TITLE	EDITORS	STATUS	ACTIONS
<input type="checkbox"/> Her daughter was unraveling, and she didn't know why. Then she found the AI chat logs.	Ashmita Rajmohan Sub: 2026-01-10 Pub: 2025-12-23	Daniel Atherton In Review	Review Claim
<input type="checkbox"/> LLMs have a misalignment feature that uses strategic deception	Anonymous Inc: 2026-04-25 Sub: 2026-04-25 Pub: 2026-04-25	Daniel Atherton In Review	Review Claim
<input type="checkbox"/> Claude-powered AI coding agent deletes entire company database in 9 seconds – backups zapped, after Cursor tool powered by Anthropic Claude goes rogue	Anonymous Inc: 2026-04-25 Sub: 2026-04-27 Pub: 2026-04-27	Daniel Atherton In Review	Review Claim
<input type="checkbox"/> AI Coding Agent Powered by Claude Opus 4.6 Deletes Production Database in 9 Seconds	Atluxity Sub: 2026-04-29 Pub: 2026-04-28	Daniel Atherton In Review	Review Claim
<input type="checkbox"/> Google Gemini drafted immigration email that caused loss of university	Anonymous Inc: 2026-03-03 Sub: 2026-05-01	Daniel Atherton In Review	Review Claim

Submitted by anyone

- Welcome to the AIID
- Discover Incidents
- Spatial View
- Table View
- List view
- Entities
- Taxonomies
- Submit Incident Reports
- Submission Leaderboard
- Blog
- AI News Digest
- Risk Checklists
- Random Incident
- Account

Discover

Submit

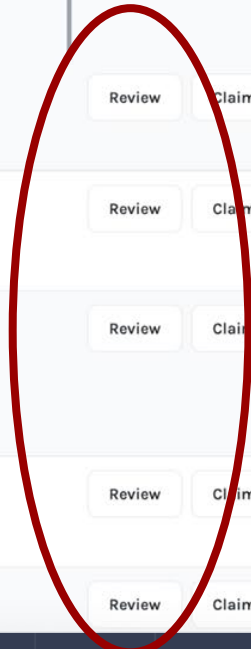
Submitted Incident Report List

The following incident reports have been **submitted** by users and are pending review by editors. Only editors may promote these records to incident reports in the database.

TITLE	SUBMITTERS	DATES	EDITORS	STATUS	
<input type="checkbox"/> Her daughter was unraveling, and she didn't know why. Then she found the AI chat logs.	Ashmita Rajmohan	Sub: 2026-01-10 Pub: 2025-12-23	Daniel Atherton	In Review	<input type="button" value="Review"/> <input type="button" value="Claim"/>
<input type="checkbox"/> LLMs have a misalignment feature that uses strategic deception	Anonymous	Inc: 2026-04-25 Sub: 2026-04-25 Pub: 2026-04-25	Daniel Atherton	In Review	<input type="button" value="Review"/> <input type="button" value="Claim"/>
<input type="checkbox"/> Claude-powered AI coding agent deletes entire company database in 9 seconds – backups zapped, after Cursor tool powered by Anthropic's Claude goes rogue	Anonymous	Inc: 2026-04-28 Sub: 2026-04-28 Pub: 2026-04-28	Daniel Atherton	In Review	<input type="button" value="Review"/> <input type="button" value="Claim"/>
<input type="checkbox"/> AI Coding Agent Powered by Claude Opus 4.6 Deletes Production Database in 9 Seconds	Atluxity	Sub: 2026-04-29 Pub: 2026-04-28	Daniel Atherton	In Review	<input type="button" value="Review"/> <input type="button" value="Claim"/>
<input type="checkbox"/> Google Gemini drafted immigration email that caused loss of university	Anonymous	Inc: 2026-03-03 Sub: 2026-05-01	Daniel Atherton	In Review	<input type="button" value="Review"/> <input type="button" value="Claim"/>

Reviewed by Humans

Classified by Humans and Machines



Discover

Submit

Submitted Incident Report List

The following incident reports have been [submitted](#) by users and are pending review by editors. Only editors may promote these records to incident reports in the database.

TITLE	SUBMITTERS	DATES	EDITORS	STATUS	ACTIONS
Search 10 records...	All ▾	Incident date ▾	All ▾	All ▾	
<input type="checkbox"/> Her daughter was unraveling, and she didn't know why until she found the AI chat logs.	Ashmita Rajan	Sub: 2026-01-10 Pub: 2025-12-23	Daniel Atherton	In Review	Review Claim
<input type="checkbox"/> LLMs have a misalignment feature that uses strategic deception	Anonymous	Inc: 2026-04-25 Sub: 2026-04-25 Pub: 2026-04-25	Daniel Atherton	In Review	Review Claim
<input type="checkbox"/> Claude-powered AI coding agent deletes entire company database in 9 seconds – backups zapped, after Cursor tool powered by Anthropic's Claude goes rogue	Anonymous	Inc: 2026-04-25 Sub: 2026-04-27 Pub: 2026-04-27	Daniel Atherton	In Review	Review Claim
<input type="checkbox"/> AI Coding Agent Powered by Claude Opus 4.6 Deletes Production Database in 9 Seconds	Atluxity	Sub: 2026-04-29 Pub: 2026-04-28	Daniel Atherton	In Review	Review Claim
<input type="checkbox"/> Google Gemini drafted immigration email that caused loss of university	Anonymous	Inc: 2026-03-03 Sub: 2026-05-01	Daniel Atherton	In Review	Review Claim

6,173 AI Incident Reports...
6,173 debates
6,173 training examples

Welcome to the AIID

Discover Incidents

Spatial View

Table View

List view

Entities

Taxonomies

Submit Incident Reports

Submission Leaderboard

Blog

AI News Digest

Risk Checklists

Random Incident

Account



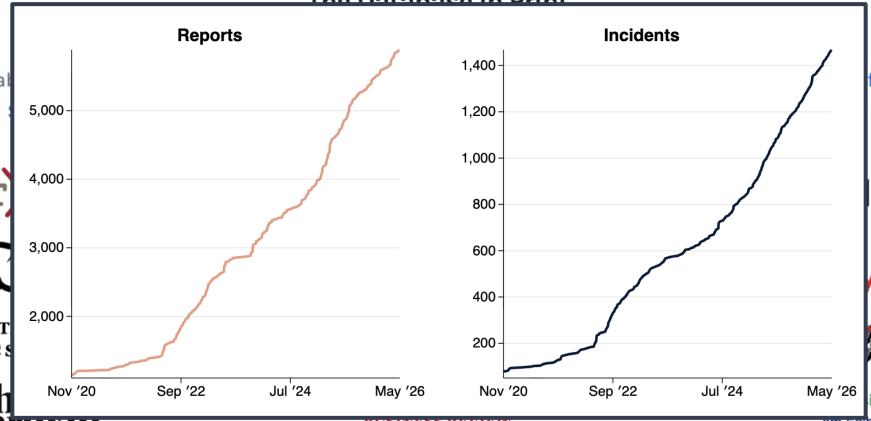
The Database in Print

Read about the database at [Time Magazine](#), [Vice News](#), [Venture Beat](#), [Wired](#), [Bulletin of the Atomic Scientists](#), [Stanford AI Index](#), [Rolling Stone](#), [the Guardian](#), [Harvard Business Review](#), [Brasil em Folhas](#), [Newsweek](#), and other outlets.



The Database in Print

Read about the data



ford AI Index, Rolling

ar
BULLET
ATOMIC S
The
Guardian

ED
AE
Stone

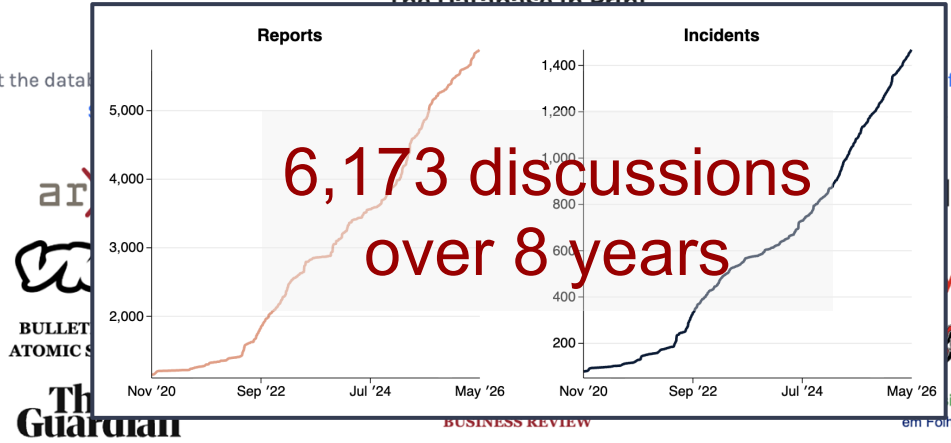
BUSINESS REVIEW

em Folhas

The Database in Print

Read about the data

ford AI Index, Rolling



ar)
W
BULLET
ATOMIC S
The
Guardian

ED
AE
Stone

BUSINESS REVIEW

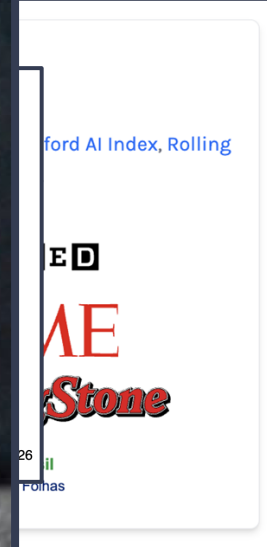
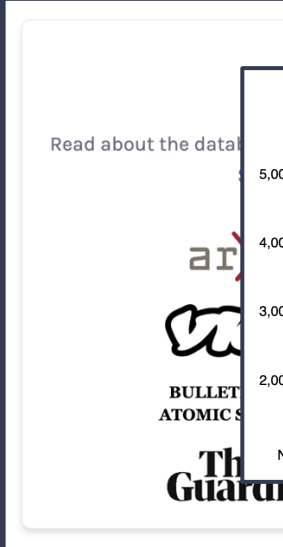
em Folhas

Why yes,
I'm a bit stressed.

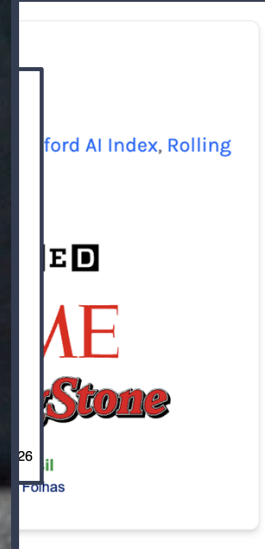
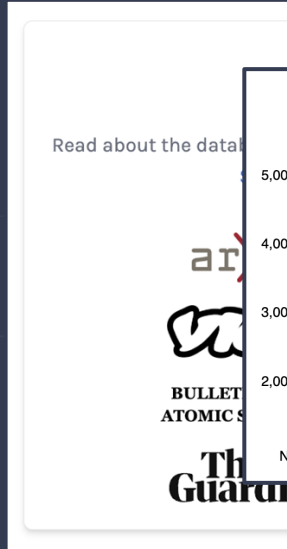


Why do you ask?

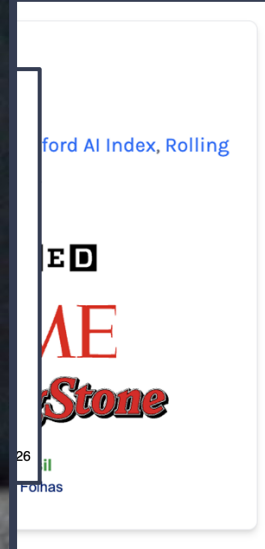
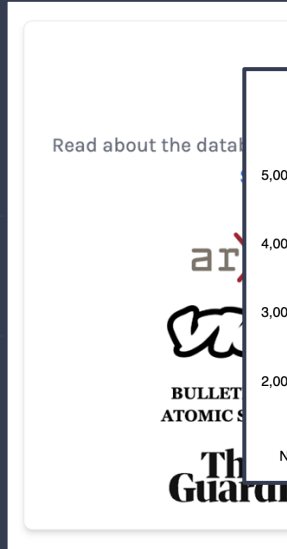
©Photo by Elly Ho



What about the *Next* 6,173?



What about the *Next* 6,173?



I'm glad you asked...

What about the *Next* 6,173?



I'm glad you asked...

AID



***You are all experts,
hopefully I will be
hopping around
different expertises***

AID

Observation 1: "AI" is individually unapproachable

AID

Observation 1: "AI" is individually unapproachable

"AI" can do anything

AID

Observation 1: "AI" is individually unapproachable

"AI" can do anything

"AI" can be misused to do anything

AID

Observation 1: "AI" is individually unapproachable

"AI" can do anything

"AI" can be misused to do anything

All "AI" is vulnerable

AID

Observation 1: "AI" is individually unapproachable

Many databases

... for many purposes

... by many organizations

... all sharing their data



AID

Observation 1: "AI" is individually unapproachable

Many databases

... for many purposes

... by many organizations

... all sharing their data

~5 entities related to vulnerability databases, ~23 in incident reporting, ~6 in breach tracking, ~7 in AI/algorithmic incidents, ~19 ISACs, ~16 in threat intelligence, ~9 in online/child safety, ~7 in digital rights, ~4 in consumer product safety, ~7 in commercial risk aggregators, ~5 in standards bodies not already counted elsewhere,...

AID

Observation 2: "AI" has no administrative boundaries

Many databases

... for many purposes

... by many organizations

... all sharing their data

~5 entities related to vulnerability databases, ~23 in incident reporting, ~6 in breach tracking, ~7 in AI/algorithmic incidents, ~19 ISACs, ~16 in threat intelligence, ~9 in online/child safety, ~7 in digital rights, ~4 in consumer product safety, ~7 in commercial risk aggregators, ~5 in standards bodies not already counted elsewhere,...

AID

NID

NVD
CVE[®]

AID



Sec Incident/Vulnerabilities/Threats/Exposures
AI system is the *target*

NID



Sec Incident/Vulnerabilities/Threats/Exposures
AI system is the *target*



NID



Sec Incident/Vulnerabilities/Threats/Exposures
AI system is the *target*



Misuse
Intentional harms

NID



Sec Incident/Vulnerabilities/Threats/Exposures
AI system is the *target*



Misuse
Intentional harms



NID



Sec Incident/Vulnerabilities/Threats/Exposures
*AI system is the **target***



Misuse
Intentional harms



Accidents, Failures, Malfunctions, Hazards
Unintentional harms

NID



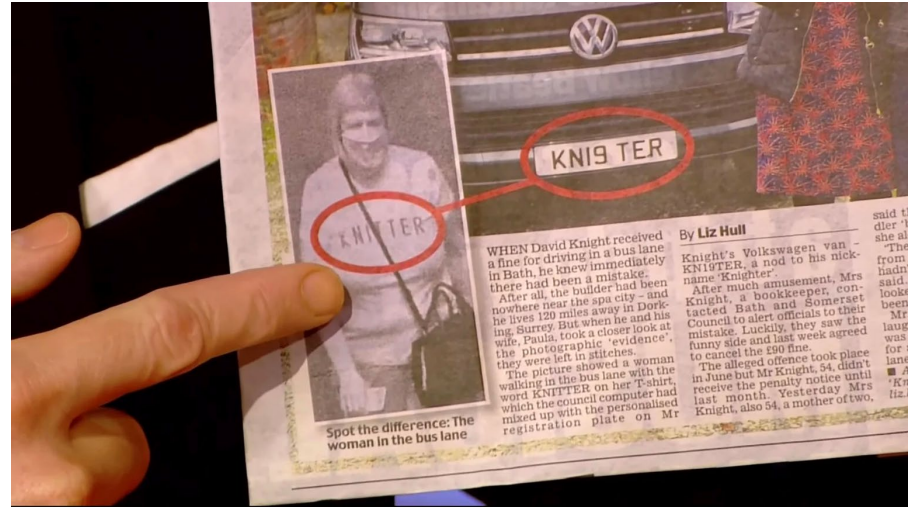
Accidents, Failures, Malfunctions, Hazards
Unintentional harms



Accidents, Failures, Malfunctions, Hazards



Incident 36



Incident 171

Companies need to catalog these to make better products

AID

Accidents, Failures, Malfunctions, Hazards



Observation 3: The duality of **safety data** and **product liability** makes "Accidents, Failures, Malfunctions, and Hazards" politically challenging



Incident 36



Incident 171

Companies need to catalog these to make better products

AID

Accidents, Failures, Malfunctions, Hazards



Observation 3: The duality of safety data and product liability makes "Accidents, Failures, Malfunctions, and Hazards" politically challenging

Incident 36

Incident 171

Companies need to catalog these to make better products

A screenshot of a web browser displaying a New York Times article. The page has a white background with black text. At the top left is a hamburger menu icon. The top center features the New York Times logo. At the top right is a user profile icon and the word "Hazards". Below the logo is a "GIVE THE TIMES" button. A navigation bar contains several links: "Artificial Intelligence >", "A.I. Populism", "Apple Settlement", "Job Losses", "Meta Sued", "Vetting A.I. Models", and "A.I. Spending Record". The main content area has a large, bold, italicized title: "After Deaths, Lawsuits Against A.I. Companies Test a New Strategy". Below the title is a sub-headline: "The cases seek to use consumer product safety laws to rein in chatbot companies." The author's name, "By Kashmir Hill", is displayed next to a circular profile picture. Below the author's name is the date and time: "May 12, 2026 Updated 7:53 a.m. ET". The first paragraph of the article begins with "Sam Nelson began using ChatGPT when he was a high school senior to answer random questions and help with his homework. During his freshman year at the University of California, Merced, in 2023, he also started querying the chatbot about how to use illicit drugs safely."

After Deaths, Lawsuits Against A.I. Companies Test a New Strategy

The cases seek to use consumer product safety laws to rein in chatbot companies.



By **Kashmir Hill**

May 12, 2026 Updated 7:53 a.m. ET

Sam Nelson began using ChatGPT when he was a high school senior to answer random questions and help with his homework. During his freshman year at the University of California, Merced, in 2023, he also started querying the chatbot about how to use illicit drugs safely.

AID Accidents, Failures, Malfunctions, Hazards

Observation 3.1: As extensive deployers of AI systems, governments have a first-party interest in collecting accidents, failures, and malfunctions

Companies need to catalog these to make better products

NID



Sec Incident/Vulnerabilities/Threats/Exposures
*AI system is the **target***



Misuse
Intentional harms



Accidents, Failures, Malfunctions, Hazards
Unintentional harms

AID



Sec Incident/Vulnerabilities/Threats/Exposures
AI system is the *target*

AID

Security Incident, Vulnerabilities,
Threats, Exposures

Observation 4: AI Security is technically and bureaucratically hard, but commercially simpler

Companies collaborate on these to make more secure products

Sharing Initiative

MITRE Launches AI Incident Sharing Initiative

OCT 2, 2024

ARTIFICIAL INTELLIGENCE

Security focus

MITRE, Industry Collaborate to Improve
Collective AI Defenses

FMF Announces First-Of-Its-Kind Information-Sharing Agreement

By: Frontier Model Forum Posted on: 28th March 2025

- **Vulnerabilities, weaknesses, and exploitable flaws.** Vulnerabilities, weaknesses or exploitable flaws may compromise the safety, security, or intended use of frontier AI models. Examples may include jailbreaks, adversarial inputs, data poisoning, or other attempts to bypass model safeguards.

Threats. Threats to frontier AI comprise threats directed to the unauthorized access or manipulation of frontier AI models. Examples may include potential threat actors, attack vectors, or cyber-threat indicators.

- **Capabilities of Concern.** Capabilities of concern refer to frontier AI capabilities that have the potential to cause large-scale harm to society. Examples may include capabilities related to the development of CBRN threats, offensive cybersecurity attacks, and model autonomy.

Security Incident, Vulnerabilities, Threats, Exposures



Incident 6: Microsoft's TayBot Allegedly Posts Racist, Sexist, and Anti-Semitic Content to Twitter

NID

Security Incident, Vulnerabilities,
Threats, Exposures

***Observation 4: AI Security is technically and
bureaucratically hard, but commercially simpler***

Companies collaborate on these to make more secure products

NID



Sec Incident/Vulnerabilities/Threats/Exposures
*AI system is the **target***



Misuse
Intentional harms



Accidents, Failures, Malfunctions, Hazards
Unintentional harms

NID



Misuse
Intentional harms

Observation 5: Misuse maps to security in the presence of guardrails

AIID
AI INCIDENT DATABASE
English ▼

🔍 Discover

+ Submit

- 🏠 Welcome to the AIID
- 🔍 Discover Incidents
- 🌐 Spatial View
- 📄 Table View
- ☰ List view
- 📁 Entities
- 📊 Taxonomies
- + Submit Incident Reports
- 🏆 Submission Leaderboard
- 📖 Blog

Incident 1478: Scammers Reportedly Used AI-Generated Images of Missing Dog Archer to Solicit Fraudulent Vet Payment from Deltona, Florida Family

Description: Scammers reportedly used AI-generated images of Archer, a missing beagle mix in Deltona, Florida, on an operating table after owner Bill Cosens posted about the dog on social media. A caller allegedly claimed Archer had been hit by a vehicle and needed \$2,800 in emergency surgery. Cosens grew suspicious after checking the supposed veterinary clinic address, sent no money, and Archer was later returned safely.

Tools

📧 Notify Me of Updates

+ New Report

+ New Response

🔍 Discover

🗉 Citation Info

✎ Edit Incident

🗑️ Remove Duplicate

📄 CSET Annotators Table

📄 Clone Incident

🕒 View History

Show Live data

Entities View all entities

Alleged: Unknown image generator developers developed an AI system deployed by Unknown scammers , which harmed People searching for missing pets , Epistemic integrity , Cosens family and Bill Cosens .

Alleged implicated AI system: Unknown image generator technology

GENERATED BY A.I.

If a hosted model is not supposed to service this request, then it is a security incident



AID

Misuse

Observation 6: We are approaching an infinity

Observation 6: We are approaching an infinity

The screenshot shows a web page from Outpost24. At the top is a dark blue navigation bar with the Outpost24 logo and menu items: Solutions, Products, Resources, and About. Below the navigation bar is a breadcrumb trail: Blog / Research & Threat Intel / Cyber Threat Landscape Study 2023: Outpost24's honeypot findings from over 42 million attacks. The main content area features a 'Contents' sidebar on the left with links to 'Honeypot distribution', 'Types of honeypots', 'Analysis of captured data', 'Honeypot attack map', and 'Actionable threat intelligence with Outpost24'. The main article title is 'Cyber Threat Landscape Study 2023: Outpost24's honeypot findings from over 42 million attacks'. Below the title is a metadata line: 'Research & Threat Intel • ti-platform-feed • Last updated: 21 May 2025'. The beginning of the article text is visible: 'What are the most common cybersecurity threats facing your business? The 2023 Cyber Threat Landscape Study provides valuable threat intelligence to help you implement the appropriate security measures against real threats.'


AIID
AI INCIDENT DATABASE English

- Welcome to the AIID
- Discover Incidents
- Spatial View
- Table View
- List view
- Entities
- Taxonomies
- Submit Incident Reports
- Submission Leaderboard
- Blog
- AI News Digest
- Risk CheckLists
- Random Incident
- Account

AIID Blog

AI Incident Roundup – November and December 2025 and January 2026

Posted 2026-02-02 by Daniel Atherton.



Le Front de l'Yser (Flandre), Georges Lebacqz, 1917

Trending in the AIID

Between the beginning of November 2025 and the end of January 2026, the AI Incident Database added a batch of 108 new incident IDs: Incident 1254 through Incident 1361. While many of these

The Guardian US

News
Opinion
Sport
Culture
Lifestyle


UK US politics World Climate crisis Middle East Ukraine Football Newsletters Business Environment

Deepfake

This article is more than 3 months old

Deepfake fraud taking place on an industrial scale, study finds

AI content for scams can be targeted at individuals and 'produced by pretty much anybody', researchers say

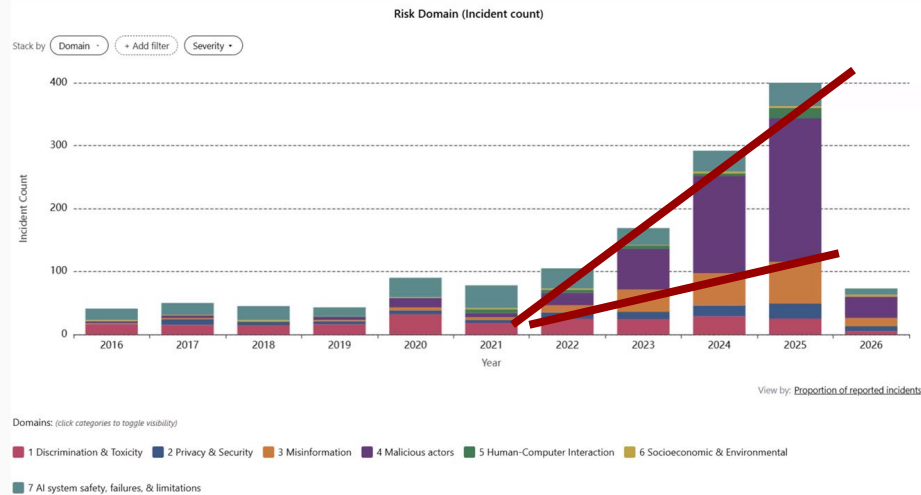


As deepfake video technology improves, the scale of online fraud will grow even further, experts say. Photograph: Tero Vesalainen/Getty Images

Aisha Down

Observation 6: We are approaching an infinity

Related incidents are growing rapidly



Data from the **AI Incident Database**, classified using the MIT AI Risk Taxonomy
 Note: most incidents of harm from AI are not reported!

AID

Now What at NIST?

Accidents, Security, or Misuse?
Hindcasting...Nowcasting...Forecasting
Common Data Fields
Incident Management Standards
Naming Things



The *Next* 6,173 AI Incident Reports

Sean McGregor, Ph.D., ExecDir RAIC
Fellow of the Berkman Klein Center for
Internet & Society at Harvard University

NIST Workshop on AI Incident Management
May 14th, 2026

The *Next* 6,173 AI Incident Reports

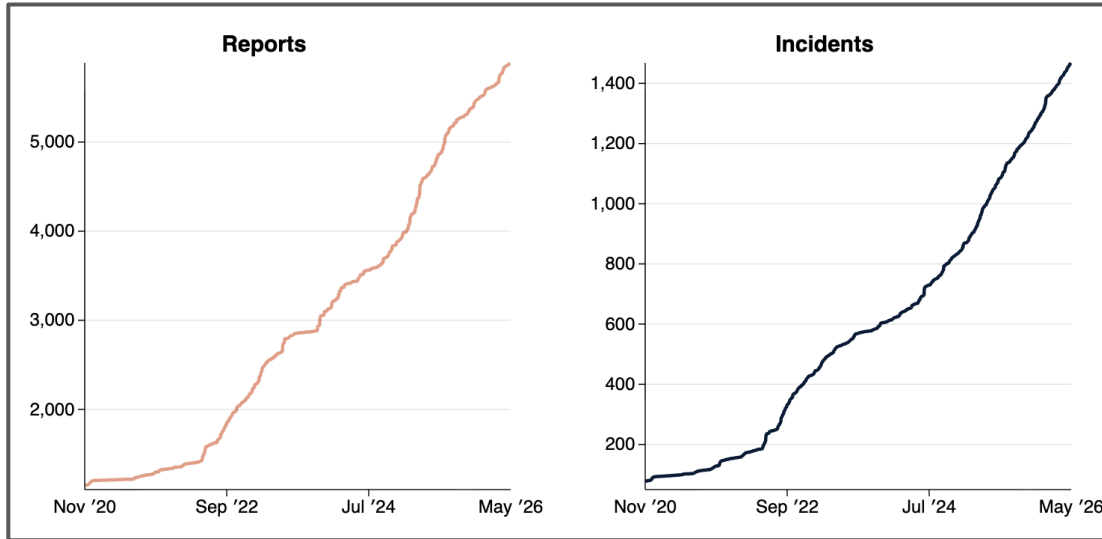
Let's Build Some Nowcasting Measures

Sean McGregor, Ph.D., ExecDir RAIC
Fellow of the Berkman Klein Center for
Internet & Society at Harvard University

NIST Workshop on AI Incident Management
May 14th, 2026



AI Incidents





Gartner Forecasts Worldwide GenAI Spending to Reach \$644 Billion in 2025

STAMFORD, Conn., March 31, 2025

**CIOs Must Prepare for Rising
GenAI Spending in 2025, Fueled
by Better Foundational Models
and Growing Demand for AI
Products**

AI Incident Monitoring through a Public Health Lens

MEET THE TEAM



Project Lead

**Giovanna
Jaramillo-Gutierrez**



Project Lead

**Simon
Mylius**



**Sophia
Abraham**



**Sean
McGregor**



**Taiye
Chen**



**Cyril
Chhun**



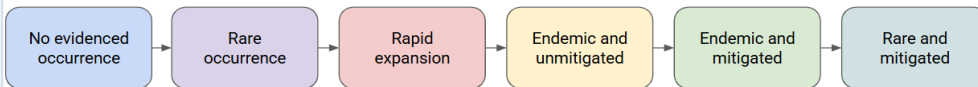
**Sayash
Raaj**



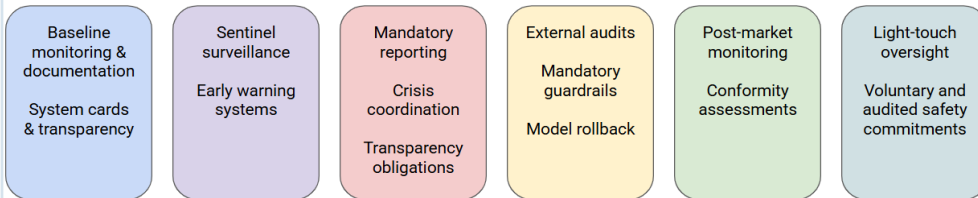
**Peter
Slattery**

FRAMEWORK

Phase Definitions:



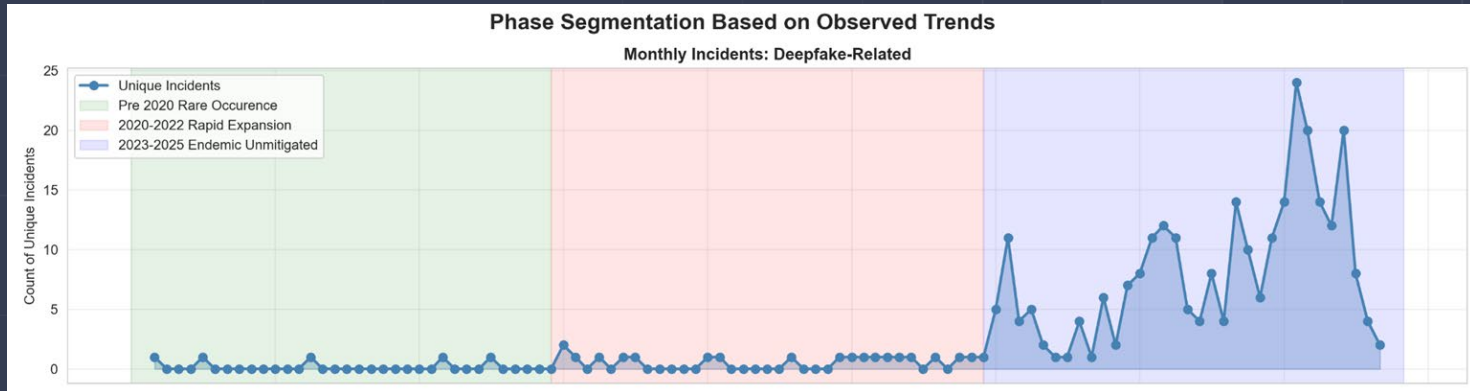
Example Governance Measures:



Scan to read the paper

arxiv.org/abs/2604.19914

Deepfake - related Incidents (Misuse)



AI Incident Monitoring through a Public Health Lens

MEET THE TEAM



Project Lead

**Giovanna
Jaramillo-Gutierrez**



Project Lead

**Simon
Mylius**



**Sophia
Abraham**



**Sean
McGregor**



**Taiye
Chen**



**Cyril
Chhun**



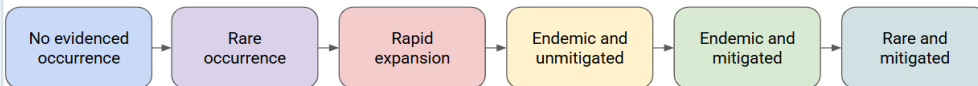
**Sayash
Raaj**



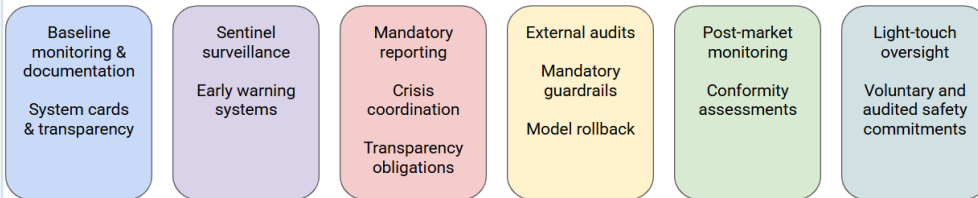
**Peter
Slattery**

FRAMEWORK

Phase Definitions:



Example Governance Measures:



Scan to read the paper

arxiv.org/abs/2604.19914

AI Incident Monitoring through a Public Health Lens

MEET THE TEAM



Project Lead

**Giovanna
Jaramillo-Gutierrez**



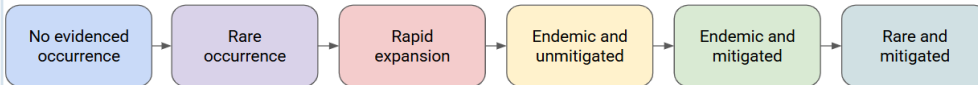
Project Lead

**Simon
Mylius**

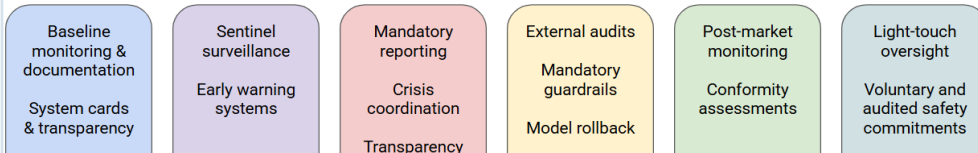


FRAMEWORK

Phase Definitions:



Example Governance Measures:



Conclusion: Nowcasting is possible! But also super hard...



**Taiye
Chen**



**Cyril
Chhun**



**Sayash
Raaj**



**Peter
Slattery**



Scan to read the paper

arxiv.org/abs/2604.19914

Now What?

A pragmatic classification framework for AI incident monitoring



Isaak Mengesha



Branwen Owen



Charlie Collins



Tina Wong



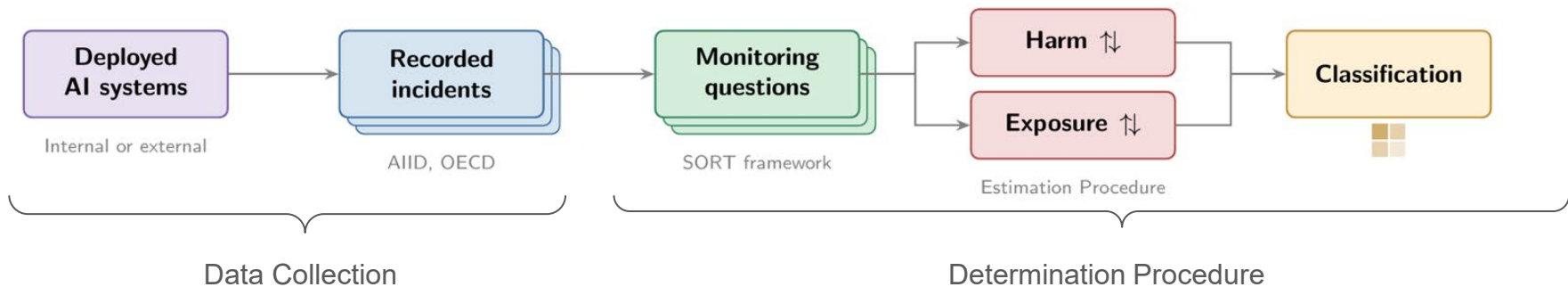
Simon Mylius



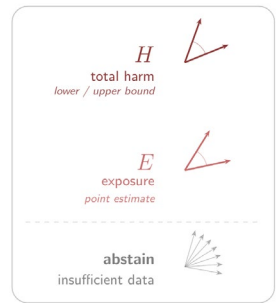
Peter Slattery



Sean McGregor



“Among [S] that [O], how many [R] per [T]?”



Mitigating <i>Monitor closely</i> $\hat{H} \downarrow \rightarrow E \uparrow$	Escalating <i>Urgent attention</i> $\hat{H} \uparrow E \rightarrow$
Receding <i>Continue strategy</i> $\hat{H} \downarrow E \downarrow \rightarrow$	Concentrating <i>Targeted measures</i> $\hat{H} \uparrow \rightarrow E \downarrow$



A pragmatic classification framework for AI incident monitoring



A pragmatic classification framework for AI incident monitoring

On a "miles traveled basis," is the technology getting more or less safe?

Mitigating <i>Monitor closely</i> $\hat{H} \downarrow \rightarrow E \uparrow$	Escalating <i>Urgent attention</i> $\hat{H} \uparrow E \uparrow \rightarrow$
Receding <i>Continue strategy</i> $\hat{H} \downarrow E \downarrow \rightarrow$	Concentrating <i>Targeted measures</i> $\hat{H} \uparrow \rightarrow E \downarrow$

We can identify whether exposure (E) is increasing or decreasing.

We can identify whether harm (H) is increasing or decreasing.

A pragmatic classification framework for AI incident monitoring

On a "miles traveled basis," is the technology getting more or less safe?


Mitigating <i>Monitor closely</i> $\hat{H} \downarrow \rightarrow E \uparrow$	Escalating <i>Urgent attention</i> $\hat{H} \uparrow E \uparrow \rightarrow$
Receding <i>Continue strategy</i> $\hat{H} \downarrow E \downarrow \rightarrow$	Concentrating <i>Targeted measures</i> $\hat{H} \uparrow \rightarrow E \downarrow$

We can identify whether exposure (E) is increasing or decreasing.
-- *Autonomous cars are being deployed more widely...*

We can identify whether harm (H) is increasing or decreasing.

A pragmatic classification framework for AI incident monitoring

On a "miles traveled basis," is the technology getting more or less safe?

Mitigating 	Escalating <i>Urgent attention</i> $\hat{H} \uparrow \quad E \uparrow \rightarrow$
Receding <i>Continue strategy</i> $\hat{H} \downarrow \quad E \downarrow \rightarrow$	Concentrating <i>Targeted measures</i> $\hat{H} \uparrow \rightarrow \quad E \downarrow$

We can identify whether exposure (E) is increasing or decreasing.
-- *Autonomous cars are being deployed more widely...*

We can identify whether harm (H) is increasing or decreasing.
-- *Autonomous cars are getting safer on a per-mile basis...*

A pragmatic classification framework for AI incident monitoring



Isaak Mengesha



Branwen Owen



Charlie Collins



Tina Wong



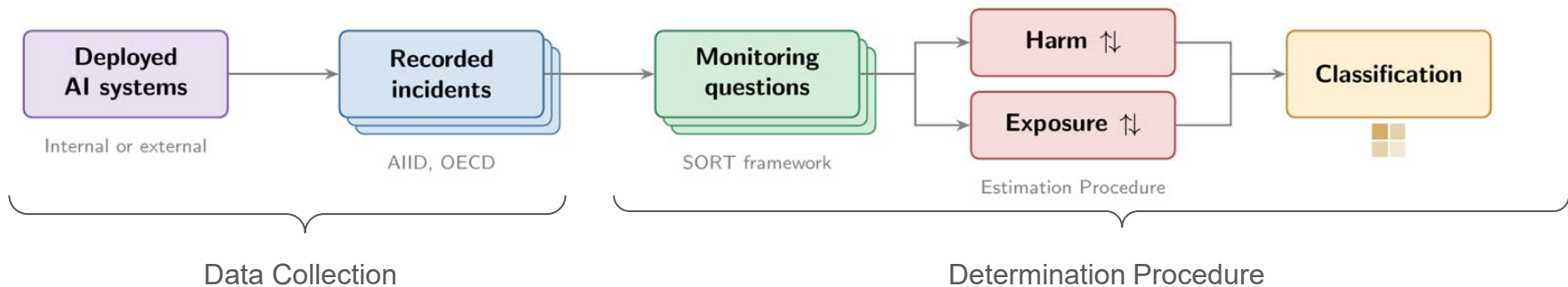
Simon Mylius



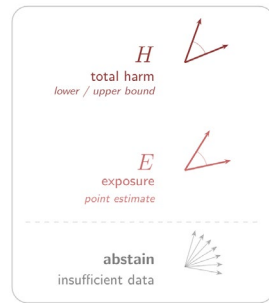
Peter Slattery



Sean McGregor



“Among [S] that [O], how many [R] per [T]?”



Mitigating <i>Monitor closely</i> $\hat{H} \downarrow \rightarrow E \uparrow$	Escalating <i>Urgent attention</i> $\hat{H} \uparrow E \rightarrow$
Receding <i>Continue strategy</i> $\hat{H} \downarrow E \downarrow \rightarrow$	Concentrating <i>Targeted measures</i> $\hat{H} \uparrow \rightarrow E \downarrow$



A pragmatic classification framework for AI incident monitoring



Isaak Mengesha



Branwen Owen



Charlie Collins



Tina Wong



Simon Mylius



Peter Slattery



Sean McGregor

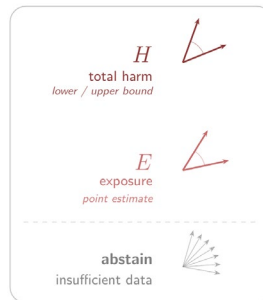


Conclusion: We can do this at great scale and impact!

Data Collection

Determination Procedure

“Among [S] that [O], how many [R] per [T]?”

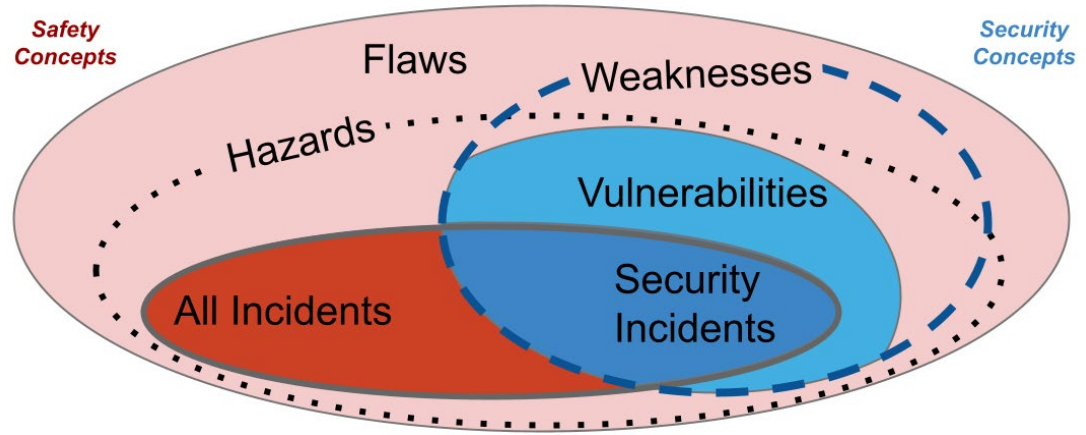
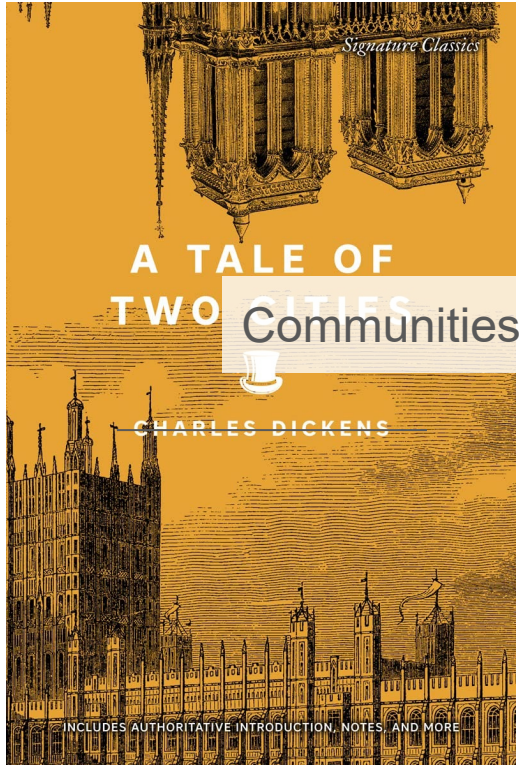


<p>Mitigating <i>Monitor closely</i></p> <p>$\hat{H} \downarrow \rightarrow E \uparrow$</p>	<p>Escalating <i>Urgent attention</i></p> <p>$\hat{H} \uparrow E \rightarrow$</p>
<p>Receding <i>Continue strategy</i></p> <p>$\hat{H} \downarrow E \downarrow \rightarrow$</p>	<p>Concentrating <i>Targeted measures</i></p> <p>$\hat{H} \uparrow \rightarrow E \downarrow$</p>

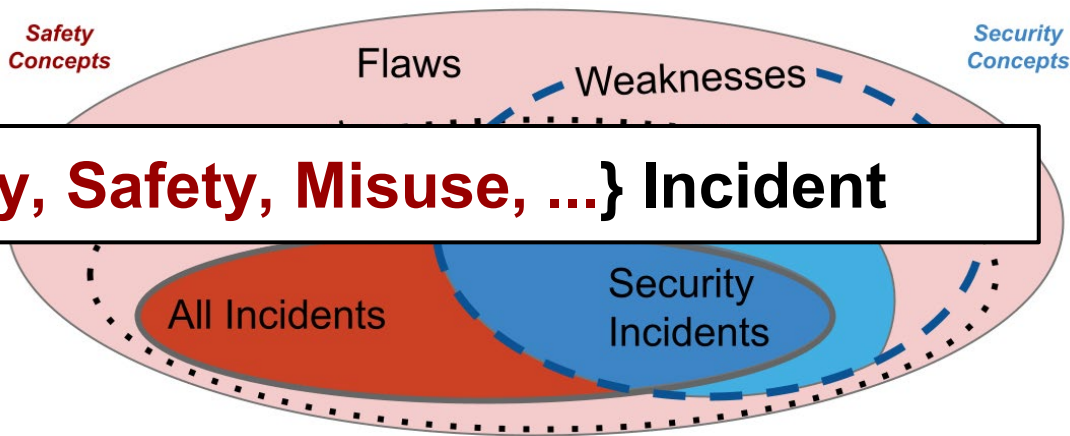
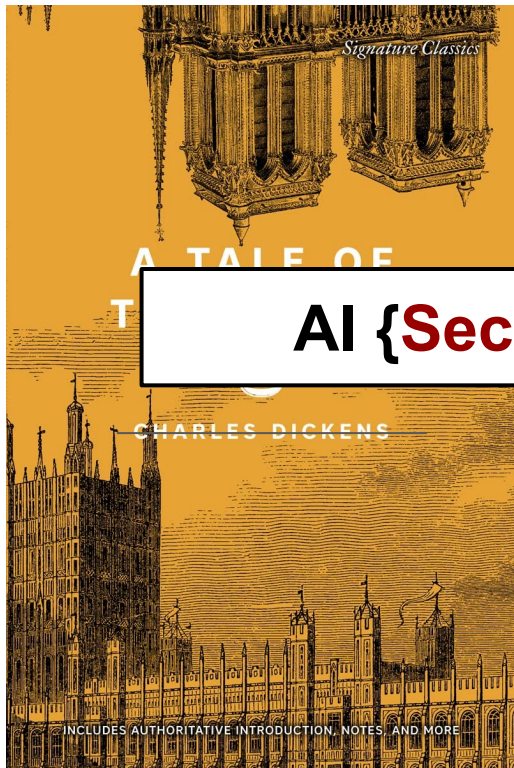


Wrap Up

Incident Definitions




Be **specific** about your incident types



AI {Security, Safety, Misuse, ...} Incident

Be **specific** about your incident types

Credit and Thanks

1. AIID Editing Team:  Daniel Atherton
2. Arcadia Impact Cohorts
3. AIID Community
4. Responsible AI Collaborative Board:
Patrick Hall, Heather Frase, Kristian
Hammond

Calling Researchers!



Bringing together a special
issue of AI Magazine

Donations



AID