

NIST WORKSHOP

# Information Sharing, Incident Reporting, and Incident Response for Frontier AI Risks

# Agenda



01 Overview

02 Information Sharing

03 Incident Reporting

04 Incident Response

05 FMF Information Sharing

06 Emerging Practices

# Information Sharing

## Overview & Description

Information sharing refers to the ongoing exchange of safety-relevant knowledge among frontier AI developers, researchers, governments, and civil society.

### Purpose

Collective learning and prevention

### Timing

Generally ongoing & continuous

### Direction

Bi-directional / multi-directional

### Key Actors

Industry, researchers, government

# Incident Reporting

## Overview & Description

Incident reporting refers to a formal notification submitted to a designated government body or authority that a qualifying AI safety or security incident has occurred.

### **Purpose**

Awareness / oversight of ecosystem

### **Timing**

Post-incident, within defined window

### **Direction**

Uni-directional, typically industry → authority

### **Key Actors**

Regulators, oversight bodies

# Incident Response

## Overview & Description

Incident response refers to the coordinated process of containing, mitigating, and recovering from an incident once one has been identified.

### **Purpose**

Containment & remediation

### **Timing**

Real-time after incident identification

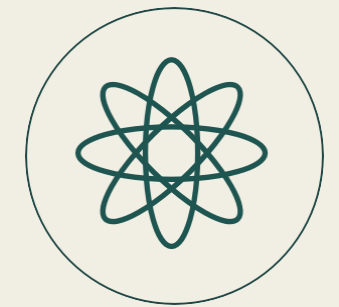
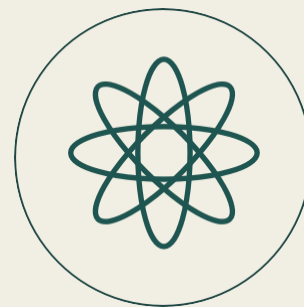
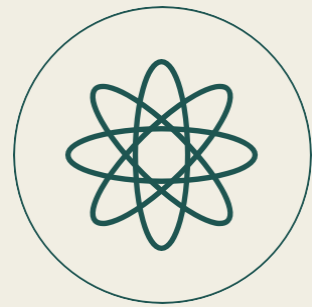
### **Direction**

Bi-directional / multi-directional

### **Key Actors**

Internal teams, affected parties

# FMF Information Sharing



## Vulnerabilities & weaknesses

Exploitable flaws that compromise the safety, security, or intended functionality of frontier AI models. Examples: jailbreaks, data poisoning.

## Threats

Threats directed to the unauthorized access or manipulation of frontier AI models. Examples: attack vectors or cyber-threat indicators

## Capabilities of Concern

Frontier AI capabilities that have the potential to cause large-scale harm to society. Examples: CBRN or advanced cyber capabilities.

The FMF has a founding mandate to establish trusted and secure information-sharing channels.  
It is **not** involved in mandatory incident reporting.

# Emerging Practices

## Institutional Design

- Define incident reporting with precision
- Build incident response procedures and capacity carefully
- Leverage existing institutional strength for incident response

# Emerging Practices

## Institutional Design

- Define incident reporting with precision
- **Build incident response procedures and capacity carefully**
- Leverage existing institutional strength for incident response

# Emerging Practices

## Institutional Design

- Define incident reporting with precision
- Build incident response procedures and capacity carefully
- **Leverage existing institutional strength for incident response**

# Emerging Practices

## Common Pitfalls

- Avoid designing information-sharing regimes that **trigger reporting obligations inadvertently**
- Avoid establishing reporting timelines that don't allow for thorough analysis.
- Avoid assuming reporting mandates substitute for response capability

# Emerging Practices

## Common Pitfalls

- Avoid designing information-sharing regimes that inadvertently trigger reporting obligations
- Avoid establishing **reporting timelines** that don't allow for **thorough analysis**.
- Avoid assuming reporting mandates substitute for response capability

# Emerging Practices

## Common Pitfalls

- Avoid designing information-sharing regimes that inadvertently trigger reporting obligations
- Avoid establishing reporting timelines that don't allow for thorough analysis.
- Avoid assuming **reporting mandates** substitute for **response capability**

# Thank you



01 Overview

02 Information Sharing

03 Incident Reporting

04 Incident Response

05 FMF Information Sharing

06 Emerging Practices