# Japan AI Safety Institute (J-AISI)

May 1st, 2025

AISI  Japan
AI Safety Institute

# Background and Overview of J-AISI

# Establishment of AISI in Japan

Following the **Hiroshima AI Process** and the UK-hosted **AI Safety Summit**, the Japan AI Safety Institute (J-AISI) was established in the IPA in Feb. 2024.

| May 2023 | November 2023 | December 2023 | February 2024 |
|---|---|---|---|
| **Agreed to the Hiroshima AI Process** "International Guiding Principles" and "International Code of Conduct | **AI Safety Summit hosted by the U.K.** | **Agreement on "Hiroshima AI Process Comprehensive Policy Framework"**  <br><br> **Prime Minister Kishida** (at the time) **announced Establishment of J-AISI** | **Japan AI Safety Institute (J-AISI) was established** |

Hiroshima AI Process
AI Safety Summit 2023
AI Strategy Council

# Integrated Innovation Strategy 2024

In the "Integrated Innovation Strategy 2024",

J-AISI is defined as **the central institution for AI Safety in Japan**.

♦ The Integrated Innovation Strategy 2024 is the fourth annual strategy that is positioned as the implementation plan for the 6th Science, Technology, and Innovation Basic Plan by the Cabinet Office.

## Three strengthening measures of the Integrated Innovation Strategy 2024

1. **Integrated strategy for key technologies**

2. **Strengthening collaboration from a global perspective**

3. **Enhancing competitiveness and ensuring safety and security in AI field**

① **AI innovation and AI accelerated innovation** (Strengthening R&D capabilities, promoting the use of AI, upgrading infrastructure, etc.)

② **Ensuring AI safety and security** (Governance, safety considerations, countermeasures against false information and misinformation, intellectual property, etc.)

③ **Promoting international cooperation and collaboration** (International cooperation based on the outcomes of the Hiroshima AI Process, etc.)

# Role and Scope of J-AISI

J-AISI's role is to support public and private sector initiatives to promote the safe and secure use of AI.

## Role

♦ Primarily plays three roles.

**Support the government**

- Investigating AI safety, examination of evaluation methods, and creating standard.

**Hub of AI Safety in Japan**

- Collecting the latest industry-academia initiatives.
- Promoting collaboration among related entities.
- Collaborating with international AI safety institutions.

**Collaboration with AI Safety-related organizations**

- Collaborate with national research institutes.
- Promote partnerships

Building a framework that enables AI developers and users to **correctly recognize AI-related risks**

**+**

Building a framework that enables **the implementation of necessary measures**, such as ensuring governance, **throughout the entire lifecycle**

↔ **Domestic & international related organizations**

▼

**Achieving a framework that balances "Promotion of Innovation" and "risk mitigation throughout the lifecycle."**

## Scope

♦ Setting the scope flexibly, while considering global trends regarding AI-related issues.

| Social Impact | Governance | AI System | Contents | Data |
|---|---|---|---|---|

# Initiatives for ensuring AI safety

J-AISI undertakes **safety evaluations**, **implementation methods**, as well as **international collaboration**.

1. **Research and evaluation of safety standards and criteria.**
   - Investigation of safety standards, check tools, disinformation countermeasures, and AI and cybersecurity.
   - Development of safety standards and guidelines.
   - Consideration of a testbed environment for AI related to the above.
2. **Study of methodologies for implementing safety evaluations.**
3. **International collaboration** with relevant organizations in other countries. (e.g. AI Safety Institute in the U.K. and the U.S.)

# Organizational Structure

AISI

# Related Government organization and agencies

AISI is a government-related organization in which 12 ministries and agencies, along with 5 related organizations, participate cross-sectionally. The secretariat is set within the IPA, under the jurisdiction of the METI* and the Digital Agency.
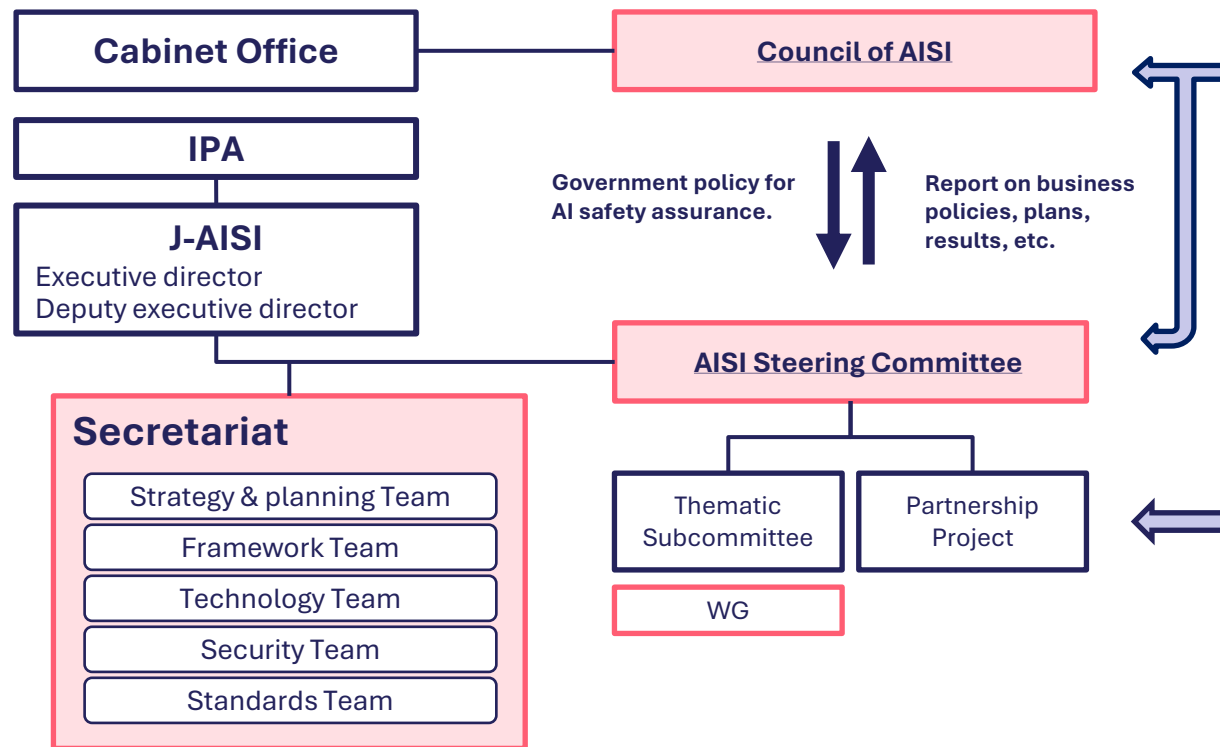
*METI: Ministry of Economy, Trade and Industry

| Cabinet Office | Cabinet Secretariat (NSS,NISC) | National Police Agency | Digital Agency | Ministry of Internal Affairs and Communications | Ministry of Foreign Affairs | Ministry of Education, Culture, Sports, Science and Technology | Ministry of Health, Labour and Welfare | Ministry of Agriculture, Forestry and Fisheries | Ministry of Economy, Trade and Industry | Ministry of Land, Infra, Transport and Tourism | Ministry of Defense |
|---|---|---|---|---|---|---|---|---|---|---|---|

12 ministries

5 organizations

| National Institute of Information and Communications Technology （NICT） | RIKEN | National Institute of Informatics （NII） | National Institute of Advanced Industrial Science and Technology （AIST） |
|---|---|---|---|

**Information-technology Promotion Agency （IPA）**

**J-AISI Secretariat**

8

# J-AISI Structure

Government policies reviewed by the AISI Liaison Meeting, led by the Cabinet Office. Project policies assessed by the AISI Steering Committee, chaired by the AISI Director.



**Relevant Ministries and Agencies:**
- Cabinet Office (Secretariat of Science, Technology and Innovation Policy)
- National Security Secretariat
- National Center of Incident readiness and Strategy for Cybersecurity
- National Police Agency
- Digital Agency
- Ministry of Internal Affairs and Communications
- Ministry of Foreign Affairs
- Ministry of Education, Culture, Sports, Science and Technology
- Ministry of Health, Labour and Welfare
- Ministry of Agriculture, Forestry and Fisheries
- Ministry of Economy, Trade and Industry
- Ministry of Land, Infrastructure, Transport and Tourism
- Ministry of Defense

**Related organizations:**
- IT Promotion Agency, Japan (IPA)
- National Institute of Information and Communications Technology
- RIKEN
- National Institute of Informatics
- National Institute of Advanced Industrial Science and Technology

# Executive Team

Executive Director **Akiko Murakami**

1999: Joined IBM Japan, Research Laboratory
2016: Joined IBM Japan, Software Development Laboratory
2021: Joined Sompo Japan Insurance Inc.
          Executive Officer, CDaO (Chief Data Officer),
          General Manager of the Data-Driven Management Promotion Department [Current]
2025: Group Chief Data Officer Executive Vice President Sompo Holdings,Inc. [Current]

Deputy Executive Director/
Secretary General

**Kenji Hiramoto**

1990: Joined  NTT DATA Corporation
2008: CIO Advisor, METI
2012: Senior Advisor to the Government CIO, Cabinet Secretariat
2021: Director of Data Strategy, Digital Agency
2023: Director, IPA Digital Infrastructure Centre [Current]
2024: Deputy Director, Secretary General of J-AISI [Current]

Deputy Executive Director

**Hideyuki Teraoka**

1999: Joined the Ministry of Posts and Telecommunications
2007: Ministry of Internal Affairs and Communications,
          Telecommunications Bureau
2023: Cabinet Secretariat, Cabinet Cyber Security Center
2024: Deputy Director of J-AISI [Current]

# Secretariat

The Secretariat is composed of the following five teams and includes many seconded personnel from government and private companies.

## Strategy & planning Team

- Strategies and planning, Budget management
- PR, Human resource development
- Coordination with domestic and international organizations

## Technology Team

- Establishment of evaluation methods for AI safety
- Development of evaluation environments

## Standards Team

- Establishment of conformity assessment methods in the AI field
- Consideration of building a domestic framework for practical implementation

## Framework Team

- Consideration of an evaluation framework for AI safety
- Coordination to ensure interoperability in AI governance

## Security Team

- Research on specific attack methods on AI systems
- Consideration of a classification system for AI security incidents
- Systematization of attacks targeting AI systems

# Activities and Deliverables

AISI

# Activities and Deliverables for FY2024

| | | International | J-AISI | Government |
|---|---|---|---|---|
| | | EVENT | DELIVERABLE | |
| 2024 | Apr | | • **JP-U.S. Crosswalk1**(4/30) | • **AI Guidelines for Business** was published(4/19) |
| | May | AI Safety Summit, Korea | | |
| | Jun | G7 Summit, Italy | | • **Integrated Innovation Strategy 2024** was published(6/4) |
| | Jul | | • Japanese Translation of U.S. AI RMF(7/4) | |
| | Aug | | • **Guide to Evaluation Perspectives**(9/18) | |
| | Sep | | • **JP-U.S. Crosswalk2**(9/18) | |
| | Oct | | • **Guide to Red Teaming Methodology***(9/25) | |
| | Nov | International Network of AISIs Convening, USA | | |
| | Dec | | | |
| 2025 | Jan | | • Published **Activity Map on AI Safety**(2/7) | |
| | Feb | AI Action Summit, France | • Published **Data Quality Management Guidebook**(Draft) (2/7) <br> • Published **National Status Report on AI Safety in Japan 2024**(2/7) | • Updated on **AI Guidelines for Business** (3/28) |
| | Mar | | | |

* Red teaming involves identifying and addressing weaknesses in AI systems from an attacker's perspective to maintain or enhance AI safety

# Summary of Deliverables

We prioritize our efforts, from technical reviews to human resource development, with the AI guidelines for Business at the center.

| **Crosswalk** | **AI Guidelines for Business** | **Activity Map on AI Safety** |
|---|---|---|
| For international interoperability | Developed and updated by the MIC* and the METI* | comprehensive overview and prioritization |
| **Guide to Evaluation Perspectives** | **Guide to Red Teaming Methodology** | **Data Quality Management Guidebook** |
| Evaluation | Red Teaming | To provide qualified data for AI |
| **Multilingual/ Multicultural** | **Security Report** | **Digital Skills Standards** |
| Multilateral challenges | Knowledge of security | Human resource development |

*MIC: Ministry of Internal Affairs and Communications
*METI: Ministry of Economy, Trade and Industry

# AI Guidelines for Business

METI* and MIC* have integrated and updated the existing guidelines, and published the AI Guidelines for business (Ver 1.0).

*Updated to Ver 1.1 in March 2025.

- Clarifying the responsibilities of each stakeholder in the process of utilizing AI.
- J-AISI co-hosts the study group for "AI Guidelines for Business" with METI.

## Content intended as a guide for businesses

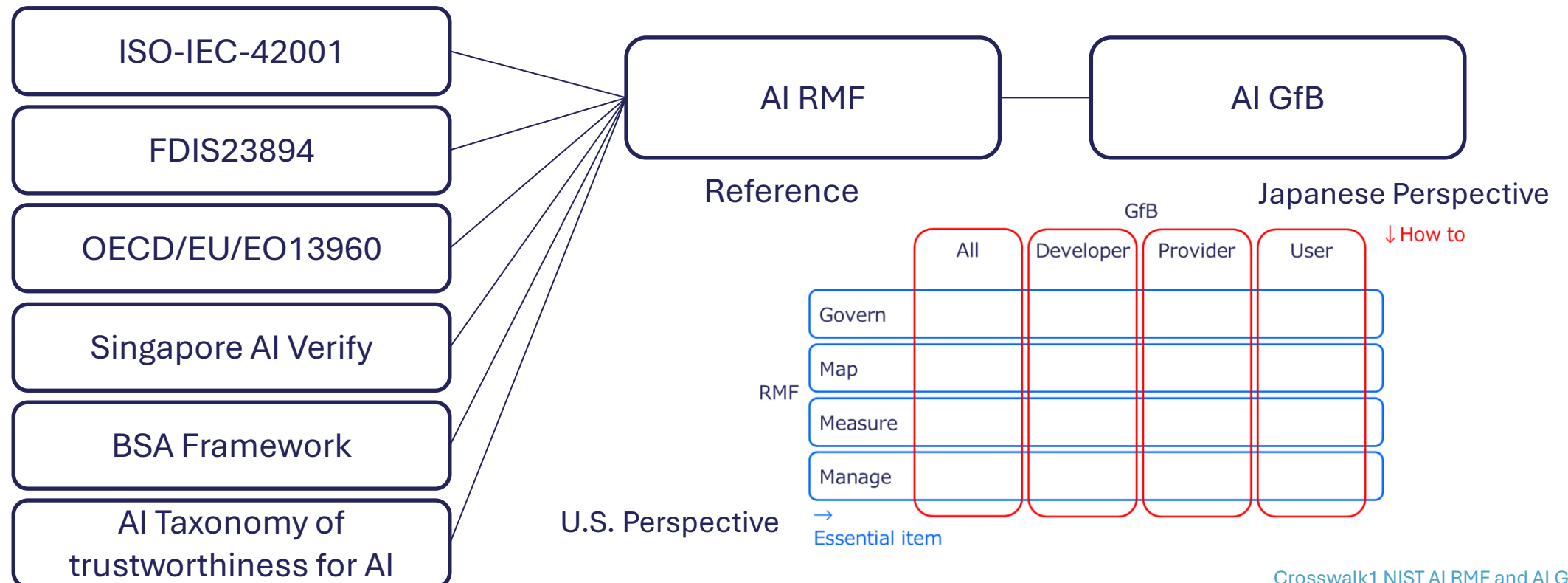### Defining the content by each lifecycle stage and entity



| AI developer | | AI provider | | AI business users |
|---|---|---|---|---|
| Data preprocessing/training | Development | Implementation in the system | Provision | Use |

Cases where an AI system is provided by an AI provider to an AI business user

AI system provision/ operational support → Proper use/ operation
- Deliver AI systems to AI business users and provide operational support such as alerts for proper use.
- Continue proper use of AI as intended by AI developers and AI providers
- Operate the AI system normally with reference to alerts from the AI provider

Data collection – Pre-processing – AI model training/ creation – AI model verification – AI system verification – AI Implementation/ coordination
- Collect data for training
- Process the collected data to apply it to AI modeling
- Create an AI model trained on the data
- Verify the usefulness of the AI model created through trials
- Implement the AI model to the system and verify it
- Integrate AI system into existing/new systems and coordinate them with other systems

Cases where an AI provider operates an AI system and an AI business user uses an AI service

AI service provision/ operation → Proper use
- Operate the AI system normally and provide AI services to AI business users
- Continue proper use as intended by AI developers and AI providers

Non-business users
- Those who use AI for non-business activities
- Those who enjoy benefits or suffer losses in some cases from AI system/service without directly using AI in their business activities

Data provider
(Not limited to cases where data is provided by a specific corporation or individual, but also includes cases where data is obtained from publicly available information on the Internet, etc.)

### Defining governance Behavioral Goals for private operators

| Category | Behavioral Goals ※ Some are further subdivided like 「3-1-1」 |
|---|---|
| 1. Environmental and risk analysis | 1-1 Understanding benefits/risks<br>1-2 Understanding social acceptance of AI<br>1-3 Understanding company's AI know-how |
| 2. Goal setting | 2-1 Setting AI governance goals |
| 3. System design | 3-1 Requiring evaluation of goal deviation and measures to minimize it<br>3-2 Improving literacy of those in charge of the AI management system<br>3-3 Enhancing AI management through cooperation between AI business actors and divisions<br>3-4 Reducing burden related to incidents involving AI Business Users and non-business users through preventive and prompt action |
| 4. Operation | 4-1 Ensuring that the operation of AI management system is explainable<br>4-2 Ensuring that the operation of each AI system is explainable<br>4-3 Considering proactive disclosure of AI governance practices |
| 5. Evaluation | 5-1 Verifying AI management system functions<br>5-2 Considering opinions of outside stakeholders |
| 6. Environment and risk reanalysis | 6-1 Reimplementing Behavioral Goals 1-1 to 1-3 at an appropriate time |

*METI: Ministry of Economy, Trade and Industry, *MIC: Ministry of Internal Affairs and Communications

AI Guidelines for Business（METI）, AI Guidelines for Business（MIC）

# Japan-U.S. Crosswalk

Confirmation of the interrelationship between the U.S. NIST AI Risk Management Framework and the Japanese AI Guidelines for Business.

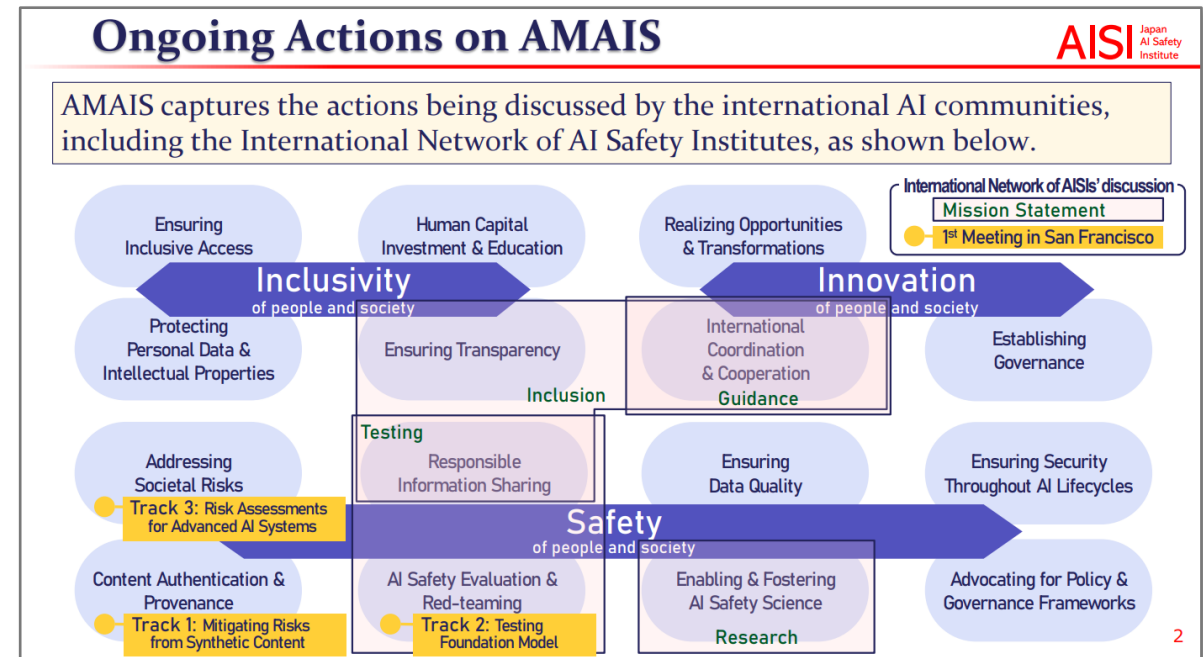♦ Using the US AI RMF as a reference, it is possible to cross-check with guidelines from other countries.



Crosswalk1 NIST AI RMF and AI GfB
Crosswalk2 NIST AI RMF and AI GfB

# AMAIS: Activity Map on AI Safety

AMAIS visualizes the overall picture of overlooked areas and correlations between activities as the AI safety efforts rapidly evolve.

♦ J-AISI has published "AMAIS: Activity Map on AI Safety" as a **discussion paper**.

♦ AISI is developing a comprehensive activity map and related terminology based on **benchmarks from key literature**.

♦ This Japan-led initiative is expected to further strengthen the foundation of international cooperation on AI safety and contribute to the realization of a sustainable and reliable AI society.

17

# Guide to Evaluation Perspectives on AI Safety

This guide presents **fundamental considerations** that can serve as a reference **when conducting AI safety evaluations.**

♦ Specifically, this document provides the following:

- **Perspectives of AI Safety evaluations, examples of risks and evaluation items**
- **Ideas on who and when it will be conducted**
- **Summary of evaluation method**

♦ This guide is a first step towards realizing safe, secure, and reliable AI, and is expected to contribute to maintaining and enhancing safety in future AI development and provision.

# Guide to Red Teaming Methodology on AI Safety

This guide serves as a resource that compiles **fundamental considerations for implementing red teaming methods** in AI safety.

- Specifically, the report provides points to keep in mind regarding the conducting structure, timing, planning, methods, and improvement plans for safety assessments.

- This guide is a first step towards realizing safe, secure, and reliable AI, and is expected to contribute to maintaining and enhancing safety in future AI development and provision.

# Data Quality Management Guidebook

The aim of this guide is to maximize the value of data/AI
and **to sustainably ensure data quality**.

- Data quality is the foundation of AI excellence, and contributes to the realization of trust AI. This guide has been organized to help us realize the trust AI society in an appropriate data quality manner and secure the data quality necessary for a data-driven society.

- The English version of this guide is the official version, while Japanese version is a summary.

# National Status Report on AI Safety in Japan 2024

> J-AISI has published the "National Status Report on AI Safety in Japan 2024", covering the **activities of J-AISI.**

- Also compiled related reports in the form of the "Fact Sheet of AI Safety in Japan 2024", a reference document that complements it.

- In the "National Status Report on AI Safety in Japan 2024", we also describe our future initiatives and aims, including collaboration between AISIs and related organizations and the private sector in Japan and overseas to respond to the rapid development of AI.

**AISI**
Japan AI Safety Institute

Please refer to the original text for accuracy

**National Status Report on AI Safety in Japan 2024**

February 5, 2025

Japan AI Safety Institute (J-AISI)

National Status Report on AI Safety in Japan 2024 - AISI Japan

# International Collaboration

# Major International meetings

AISI Japan AI Safety Institute

Actively participating in international meetings and engaging in discussions with AI-related businesses and organizations worldwide.

♦ **AISI-related collaborations**

- **Stanford University AI Symposium** (Stanford/ Apr 16, 2024)
  - Panel discussion with directors of U.S. and U.K. AISI, and parallel exchange of opinions among countries
- **AI Seoul Summit** (Seoul/ May 21-22, 2024)
  - High-level roundtable and exchange of views with the U.S., U.K., EU, Canada, Germany, etc.
  - Participation in discussions including Asian and African countries at the concurrent AI Global Forum.
- **UN Future Summit** (UN/ Sep 22, 2024)
- **UN Global Compact Leaders Summit 2024** (UN/ Sep 24, 2024)
- **AISI International Network Convening** (San Francisco/ Nov 10-11, 2024)
- **AI Action Summit** (Paris/ Feb 6-11, 2025)
- **Hiroshima AI Process Friends Group Meeting** (Tokyo/ Feb 27-28, 2025)

AI Global Forum, held concurrently with the AI Seoul Summit

UN Future Summit

**AISI** Japan AI Safety Institute

The status of AISI establishments in various countries is as follows.

*Report by CSIS in October 2024.

| | United States | United Kingdom | European Union | Japan | Singapore | South Korea | Canada |
|---|---|---|---|---|---|---|---|
| **Established** | February 2024 | November 2023 | May 2024 | February 2024 | May 2024 | May 2024 (Announced) | April 2024 (Announced) |
| **Name of Organization** | US AISI | UK AISI | EU AI Office | Japan AISI | Singapore AISI | Korea AISI | Canada AISI |
| **Housed Under** | National Institute of Standards & Technology | Department for Science, Innovation & Technology | Directorate General for Communications Networks, Content and Technology. | Information-Technology Promotion Agency | Digital Trust Centre | Electronics and Telecommunications Research Institute | |
| **Funding** (USD & Foreign Currency) | $10 million (FY24) | > $65 million/yr (>£50 million/yr 2024-2030) | $51 million (€46.5 million) (Funding period unknown) | | $7.5 million/yr (S$10 million/year) (2023-2027) | $7.2-14.4 million/yr (W10-20 billion/yr) (Tentative, starting 2025) | $36.5 million (C$50 million) (Funding period unknown) |
| **Staff** | c.20 (current core staff) | c.20 (current core staff) | c.50 (planned, AI safety unit) | c. 23 (current staff) | | Minimum 30 staff (planned, budget pending) | |
| **Public List of Functions** | US AISI Vision, Mission, and Strategic Goals | Introducing the AI Safety Institute | Tasks of the AI Office | AISI's Tasks | Initial Research Areas | | |
| **Published Research or Guidelines** | Managing Misuse Risk for Dual-Use Foundation Models | See website | | See website | Model AI Governance Framework for Generative AI | | |

| Legend | |
|---|---|
| No public statement | |
| Public Information | |

# AISI
## Japan AI Safety Institute