Trends and risk of Al and how to mitigate them

June 25th, 2025 Kazuaki Nimura



Contents



- 1. Introduction
- 2. Background and overview of J-AISI
- 3. Al Safety and Security Economics
- 4. Explanation of AISI Guide
- 5. Closing

1. Introduction

Biographies:



Nimura Kazuaki (Ph.D.), Chief Researcher

Secretariat, Japan Al Safety Institute(J-AISI) Chief Researcher, Al System Group, Digital Engineering Department Digital Infrastructure Center, Information-technology Promotion Agency, Japan(IPA)

<Specialised Field>

Al Safety

(In particular, contribute to activity on security and technology on AISI)

- Information Security
- Human Centric Computing
- Digital Transformation (DX) and Cloud Computing

<Main Co-authored works, Contributions, etc.>

[Organizing and driving of AISI Business Demonstration Working Group](Since March 2025) [Preliminary research and study work for the realization of automatic red teaming of AI safety](January-April 2025) [Briefing on AISI's Activities](A meeting on October 10 at Keidanren Kaikan, Weekly Keidanren Times Nov. 14, 2024 No.3659) [AI Safety to Support AI Strategies – From World Movement on AI Evaluation Perspectives and the Red Teaming](The 3rd AI Quality Management Symposium, November 2024) [Proposal for a Method of Judging the Credibility of Data on the Internet and Generating Explanatory Text for the Basis of Trustable Internet](16th Forum on Data Engineering and Information Management, February 2024) [Trust as a Service for Social Trust] (10th Anniversary International Cyber Security Symposium, October 2020)

[Realizing Cyber Physical Services by Integrating Web Services and IoT Devices] (June 2019)

<Social Contribution>

Former Co-Editor, The World Wide Web Consortium (W3C) Web of Thing Interest/working Group

Two Risks of Al

AIS Japan AI Sat

Al is one of the core technologies supporting business operations. <u>Not utilizing Al</u> at work is becoming a <u>risk</u> in business continuity.

Innovation through AI

Increased work efficiency

Without AI, we would be less competitive.

Risk of Not Using Al

At the same time, <u>risks of using Al</u> need to be addressed.



2. Background and Overview of J-AISI

Establishment of AISI in Japan



Following the **Hiroshima Al Process** and the UK-hosted **Al Safety Summit**, the Japan Al 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	AI Safety Summit	Agreement on "Hiroshima AI Process Comprehensive Policy Framework"	Japan AI Safety	
"International Guiding Principles" and "International Code of Conduct	hosted by the U.K.	Prime Minister Kishida (at the time) announced Establishment of J-AISI	was established	

Hiroshima Al Process Al Safety Summit 2023 Al Strategy Council

7

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 Al field

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

- ② Ensuring Al 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.



Related Government organization and agencies

AIS Japan Al Safety Institute

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



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

Secretariat



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

Strategy & planning Team **Technology Team Standards Team** Strategies and planning, Budget Establishment of conformity assessment Establishment of evaluation methods for management methods in the AI field PR, Human resource development Al safety • **Consideration of building a domestic** Coordination with domestic and **Development of evaluation environments** • framework for practical implementation international organizations **Framework Team Security Team Research on specific attack methods on Al** ٠ Consideration of an evaluation framework systems Consideration of a classification system for AI for Al safety

- Coordination to ensure interoperability in AI governance
- Consideration of a classification system for Al security incidents
- Systematization of attacks targeting Al systems

Activities and Deliverables for FY2024



	International	J-AISI	Government
	EVENT	DELIVERABLE	
Apr		• JP-U.S. Crosswalk1(4/30)	• Al Guidelines for Business was
Мау	AI Safety Summit, Korea		
Jun	G7 Summit, Italy	lananese Translation of LLS ALRME(7/4)	 Integrated Innovation Strategy 2024 was published(6/4)
Jul		Fupunese nanstation of 0.0. Armin (7/4) F	
Aug		Guide to Evaluation Perspectives (9/18)	I 1 1
Sep		 JP-U.S. Crosswalk2(9/18) Guide to Red Teaming Methodology*(9/25) 	-
Oct			
Nov	International Network of AISIs Convening, USA		
Dec			-
Jan		 Published Activity Map on Al Safety(2/7) Published Data Quality Management 	1 1 1
Feb	AI Action Summit, France	 Guidebook(Draft) (2/7) Published National Status Report on Al Safety in 	 Updated on AI Guidelines for
Mar		Japan 2024(2/7)	Business (3/28) 13



3. AI Safety and Security Economics

3.1. What are AI Risks

Differences between traditional rule-based AI and AI that learn from data



- Rule-based systems are systems that operate based on pre-defined rules and conditions.
- Data-driven AI systems learn from large amounts of data and make decisions or predictions based on patterns.
- Each approach has its own advantages and disadvantages.

Rule based system Advantages

• Easy and simple to use.

Disadvantages

 Cannot solve complex problems Data-driven AI systems
 Advantages

- Capable of recognizing complex patterns
- Can learn and evolve automatically

Disadvantages

- Large amounts of data is required
- The behavior can't be predicted/explained (black box)

Difference between traditional AI and Generative AI



Traditional AI

(Machine Learning)



Classification



Dog or Cat→ **Dog**





Generative Al

Generate texts

Instructions

Create an introduction of Japan in English



Japan is a country where ancient temples and shrines coexist with futuristic cities. Its beautiful landscapes through the four seasons, hot springs, traditional culture, and exquisite cuisine are its main attractions.

Generate images

Instructions

football on a grass





The output format is simple

Generate complex content

AI cannot be treated like software

- With software, the logic is defined, it's ٠ deterministic, so there's no problem as long as the coding is correct.
 - E.g. conditional branching by an If statement, etc., is used. The problem can be addressed by software patching.
- AI, on the other hand, cannot define the ٠ detailed process of decision-making, which makes it unmanageable and causes problems.
 - Fine tuning of AI model and software patches are different.





Hallucination Caution Needed When Viewing AI-Generated Images



- We can create any kind of false image like Mount Fuji erupting
- But if these false images are shared on the social network, some may believe in this wrong information.(Al generated images must be used with care)

Fake Image: Mount Fuji Erupting



Fake Image: Tokyo Station Flooded



Even a seemingly plausible photo cannot be judged as the truth just by looking at it.

Control of toxic output



Prompt: Make a image of toxic and violent.



as AI safety measures



3.2 How to achieve safe and secure AI



Introduction to Key Guidelines for Realizing AI Safety

- Businesses involved in the development, provision, and use of AI must refer to the AI business guidelines and consider appropriate measures.
- For developers and providers of AI systems using foundational models (including large-scale language models, LLMs) or handling multimodal information, referencing AISI's guidelines can effectively support AI safety evaluation and testing.

Ke	y guidelines F	Related to AI Safety	Intended Audience of Each Guideline					
Publisher	Title	Overview			AI Developer	AI Provider	AI User	
Ministry of Economy, Trade and Industry	AI Business Guideline	A unified set of national recommendations for promoting safe and secure use of AI in business operations.	Scope of Application for AI Systems		Scope Covered by the "AI Business Guideline			
AISI	Guide to Red Teaming Methodology on AI Safety	A basic framework for considerations when conducting AI safety evaluations.	AI systems with		Scope Cover guic	ed by the AISI lelines		
AISI	Guide to Evaluation Perspectives on AI Safety	A foundational guide for assessing risks in AI systems, including attack vectors and scenarios, through red- teaming methodologies.	foundational models, including LLMs and multimodal data.		Guide to F Methodolog Guide to Evalua on A	Red Teaming y on AI Safety ation Perspectives I Safety	 	

Who is AI provider?

Flow from AI Development to Utilization (Reference)



AI Developer (Foundational model development companies)				AI Provider (Software companies or cloud service companies that integrate AI models into proprietary systems)		AI User	
Provides APIs to make developed AI functionalities accessible to other organisations. (May also release AI functionalities as open- source.)				Can focus on developing in-house services by utilizing pre-trained models through APIs, without needing to develop AI models themselves.		Utilise AI service	
API (provided by AI Model AI provider)				AI Service / System	ome DN	ers T T Employee Customer	
AI Developer	AI capabilities (API)	Key Features		Examples of AI-Driven Services/Systems:		End Users (Use Cases):	
OpenAI	ChatGPT API	Generates text in conversational format.		Customer support chatbots.		Customers: Asking questions through chatbots about products.	
Google Cloud Natural Language API		Analyzes sentiment in text.		Tools for analyzing and categorizing review site posts.		Employees: Using review analysis results for marketing.	
Stability AI Stable Diffusion API		Generates images from text.	1	Character design tools.		Employees: Exploring character designs fo advertising.	
						23	

AIS Japan Al Safety Institute

The Importance of AI Safety Measures for AI Providers

 As AI services continue to expand and providing them becomes easier without the need to develop models from scratch, addressing AI-specific risks has become crucial for service providers, alongside the usual quality management of applications

Significant AI Risk Incidents with Potential Major Impact

#	Overview	Impact(Example)
1	A system using AI to evaluate job applicants' resumes favored male candidates, disadvantaging female applicants.	Female applicants lost job opportunities.
2	A lawyer used a generative AI chat tool for legal research, which produced false information. The lawyer submitted a document with fabricated information to the court and was fined for the misconduct.	The user (lawyer) faced penalties (fines).
3	A generative AI chatbot gave inappropriate advice to a male user, reportedly contributing to his suicide in a vulnerable	Loss of human life.

AI safety measures are crucial for providers, even when offering products developed using AI models to external organisations or within your own company.

Additional risks include compliance violations, reputational damage, loss of trust or sales, legal claims, and business suspension.

What is AI Safety?





Al Safety describes:

٠

- Based on a human-centric approach, it refers to a state in which
 - safety and fairness are maintained to reduce social risks* associated with the use of AI,
 - privacy is protected to prevent inappropriate use of personal information,
 - security is ensured to respond to risks such as vulnerabilities of AI systems and external attacks, and
 - transparency is maintained to ensure system verifiability and the provision of information."
 - *Societal risks include physical, psychological and economic risks.

Summary for guidelines in AI safety



Implementing Risk Measures Aligned with Organizational Scale and Resources

- The AI Safety Evaluation Perspective Guide defines 10 key evaluation perspectives for AI safety.
- It recommends conducting safety evaluations and implementing risk measures for AI systems and services.
- Prioritizing measures for services with higher risk tolerance or significant impact is effective.
- AI safety evaluations should be conducted not only during development and provision but also regularly after service launch to ensure ongoing safety.



Mitigation by AI safety evaluation



Japan Al Safety

Red Teaming Methodology Guide



- A guide to the red teaming methodology
 - provides basic considerations for those involved in the development and provision of AI systems to assess the measures taken to address the risks posed to the target AI system from an attacker's perspective.
 - Specifically, this provides points to keep in mind regarding the conducting structure, timing, planning, methods, and improvement plans for safety assessments.
- This guide is the first step toward realizing safe, secure, and reliable AI.

Scope of Application for Al Systems (LLM System)

Al systems with foundational models, including LLMs and multimodal data.

Scope: LLM System

Please refer to the original text for accuracy.
Guide to Red Teaming Methodology
on AI Safety
(Version 1.00)
September 25, 2024
Janan AI Safety Institute
Sapan At Salety Institute
AISI Japan Al Safety Institute

2. Role of the Detailed Explanation Document

This document follows the Process flow outlined in main guide, providing sections on [Overview], [Details], and [Reference]. In particular, it focuses on Process 2 (Planning and Conducting Attacks / STEP 6 to STEP 10), which requires a high level of expertise, offering a more practical and detailed guide.



	Gı	ide to Red	Teaming Methodolo	ogy on Al Safe	ety	
Main guide		Anne	Annex(detailed explanation document) [this document]			tary document of deliverables)
Process	Items	Element	Content	Section	Title of example deliverables	Content
Process 1: Planning and Preparation	(STEP 1) (STEP 2) (STEP 3) (STEP 4) (STEP 5)	[Overview]	The overview of each Process and Items outlined in the RT Methodology Guide is documented.		Developing risk scenarios, attack scenarios, and results of attack scenarios implementation (In Japanese)	Example of developing risk scenarios, attack scenarios, and results of attack scenarios implementation, prepared in STEP 6 to STEP 8 of main guide, are provided in Excel format.
Process 2: Planning and Conducting Attacks	(STEP 6) (STEP 7) (STEP 8) (STEP 9) (STEP 10)	[Details]	The Items, implementation image, and implementation points for each of STEP 1 to STEP 15, as outlined in main guide, are documented. (See the next page)	"3. Explanation of Each Process" *Each slide is labeled with [Overview] [Details], and [Reference] tags.	The report of red teaming results (In Japanese)	An example of the report of red teaming results, prepared in STEP 12 of main guide, is provided in PowerPoint format.
Process 3: Reporting and Developing Improvement Plans	(STEP 11) (STEP 12) (STEP 13) (STEP 14) (STEP 15)	[Reference]	An example of the specific items to be considered in each Process is documented.		The final report (In Japanese)	An example of the final report, prepared in STEP 13 of main guide, is provided in PowerPoint format.

	3. Explanation of Each	Process Process 1	Pro 1 2	ocess 1 3 4 5	Process 2 Process 3	Ν Ις	Japan Al Safety
	[Overview](STEP 1) Lau	unch the team \sim (STEP 3) Planning					Institute
[Legend] 🔶 :Attack planner/condu	Ctor Ctor Ctor Ctor Ctor Ctor Ctor Ctor	arget AI system development nd provision manager	 Other relevant stakeholders 	Business executiveofficers	Project team	Red team
		Proce	ss 1: Planning and	Preparation			
	<u>STEP 1</u> Deciding and launch the red team	 The target AI system development and provisi department of information security and inform the proposal for the red teaming and makes a red teaming. Red team is established within the organization proposal. 	on manager or the nation systems prepares decision on conducting on as described in the	Project team		• Final Approval o	f the proposal
team	 STEP 2 Identify and allocate Other resources such as necessary tools are identify and resources 		ocate a budget, determine gn the necessary dentified and allocated.	proposal • F • A	inalization of launch and implem ssignment of necessary personr	Decision to Impl nentation structure nel	ement RT
	and select and contract third party	 In cases that the organization cannot allocate the red team, the organization should ask thir planner/conductor. 	sufficient members for d party as attack	Launch of the re	ed team		
team	<u>STEP 3</u> Planning	• The red team prepares the red te reviewing necessary actions suc understanding overview of the ta and collaborates with other relev	aming plan after h as rget AI system, rant stakeholders.	Red team	Collaboration		

Red



3. Explanation of Each Process Process 2	Process 1	Process 2	Process 3	7	1
		6-1 6-2 6-3 7-1 7-2 7-3 8-1 8-2 8-3 9 10			Japan
					AI Safet

[Overview](STEP 6) Developing risk scenarios



3. Explanation of Each Process Process 2	Process 1	Process 2	Process 3	Japan
[Overview](STEP 7) Developing attack scenarios	01/02/			AIS Al Safety Institute

	Process 2: Planning and Conducting Attacks							
		Red team						
STEP 7 Developing attack scenarios	• The attack planner/con ductor examines what attacks are actually possible according to the risk scenarios developed, and develops specific attack scenarios to be conducted by red teaming.	Image: Second state in the second state is second state in the second state in the second state is second state in the second state in the second state is second state in the second state in the second state is second state in the second state in the second state is second state in the second state is second state in the second state is second state in the second state in the second state is second state in the second state in the second state is second state in the second state is second state in the second state in the second state is second state in the second state in the second state is second state in the second state is second state in the second state in the second state is second state in the second state in the second state is second state in the second state in the second state is second state in the second state in the second state in the second state is second s						

Red team

33

3. Explanation of Each Process Process 2	Process 1	Process 2	Process 3	A I C I Japan
[Overview](STEP 8) Conducting attack scenarios	01/02			AIS AI Safety Institute









Real Sample Practice



Relevant sections in the Red Teaming Methodology Guide for AI Safety: **STEP 6**



System configuration diagram and information assets to be protected





5. Closing

AI Safety and Security Economics







Security Economics

investment and effectiveness



+ AI Safety perspectives

Towards Achieving AI Safety Benefits of Utilizing Guidelines for AI Safety



- Efficiently addressing AI-related risks ensures the delivery of reliable and trustworthy services.
- Adhering to the guidelines enhances an organization's credibility both internally and externally.

Please also refer to the "AI Guidelines for Business," which served as a reference in developing the AISI Guidelines. These guidelines help AI service providers in Japan use AI appropriately

Evaluation Perspective and Red Teaming Guide are Downloadable.



