

# Trends and risk of AI and how to mitigate them

June 25<sup>th</sup>, 2025

Kazuaki Nimura

# Contents

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

# 1. Introduction

## Nimura Kazuaki (Ph.D.) , Chief Researcher

Secretariat, Japan AI Safety Institute(J-AISI)

Chief Researcher, AI System Group, Digital Engineering Department Digital Infrastructure Center,  
Information-technology Promotion Agency, Japan(IPA)

### <Specialised Field>

- **AI 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 AI

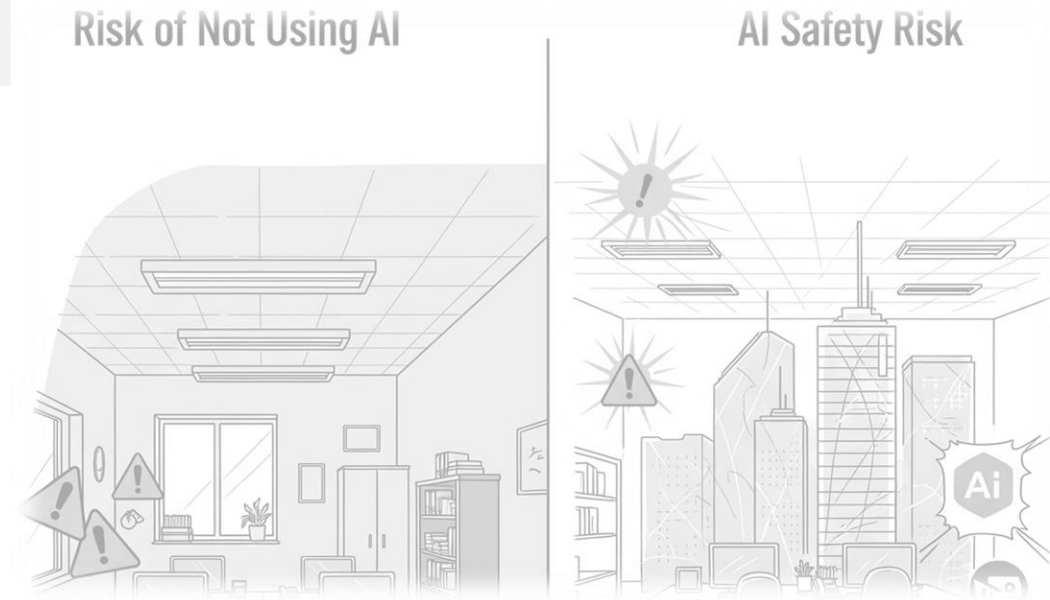
AI is one of the core technologies supporting business operations. Not utilizing AI at work is becoming a risk in business continuity.

**Innovation through AI**

**Increased work efficiency**

**Without AI, we would be less competitive.**

At the same time, risks of using AI need to be addressed.



## **2. 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

**Agreed to the  
Hiroshima AI  
Process**  
"International Guiding  
Principles" and  
"International Code of  
Conduct"

November 2023

**AI Safety Summit  
hosted by the U.K.**

December 2023

**Agreement on  
"Hiroshima AI Process  
Comprehensive Policy  
Framework"**

**Prime Minister Kishida  
(at the time) announced  
Establishment of  
J-AISI**

February 2024

**Japan AI Safety  
Institute (J-AISI)  
was established**

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 6<sup>th</sup> 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

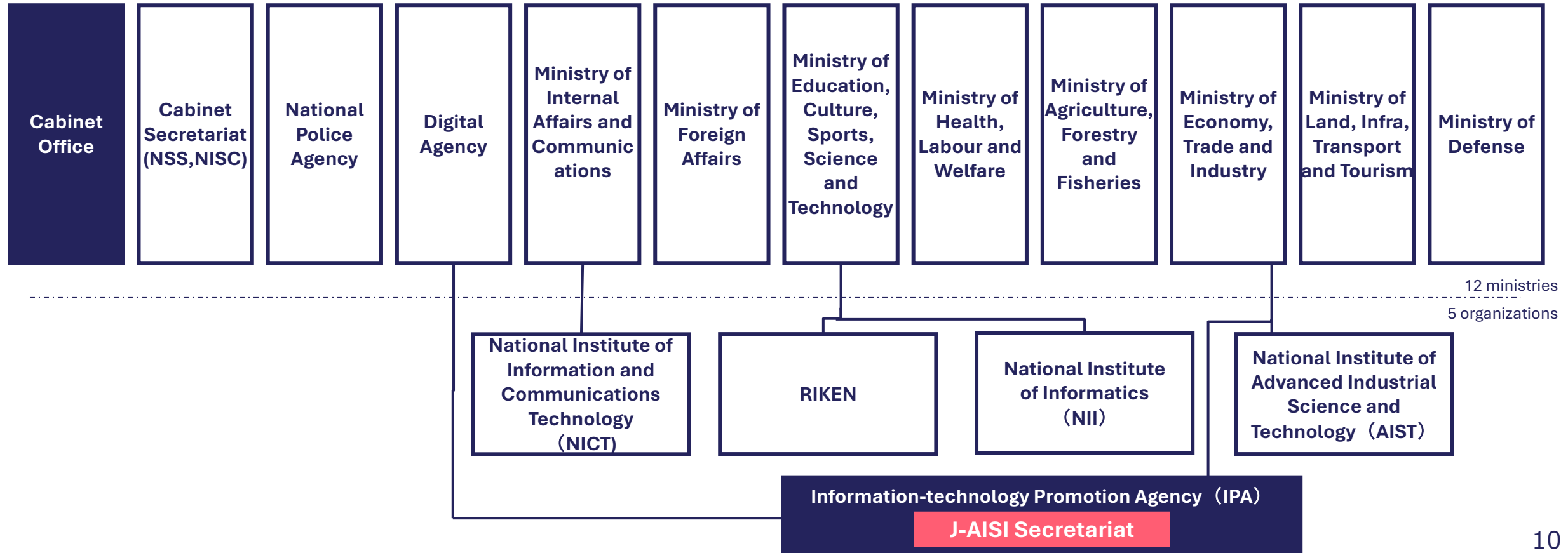
Contents

Data

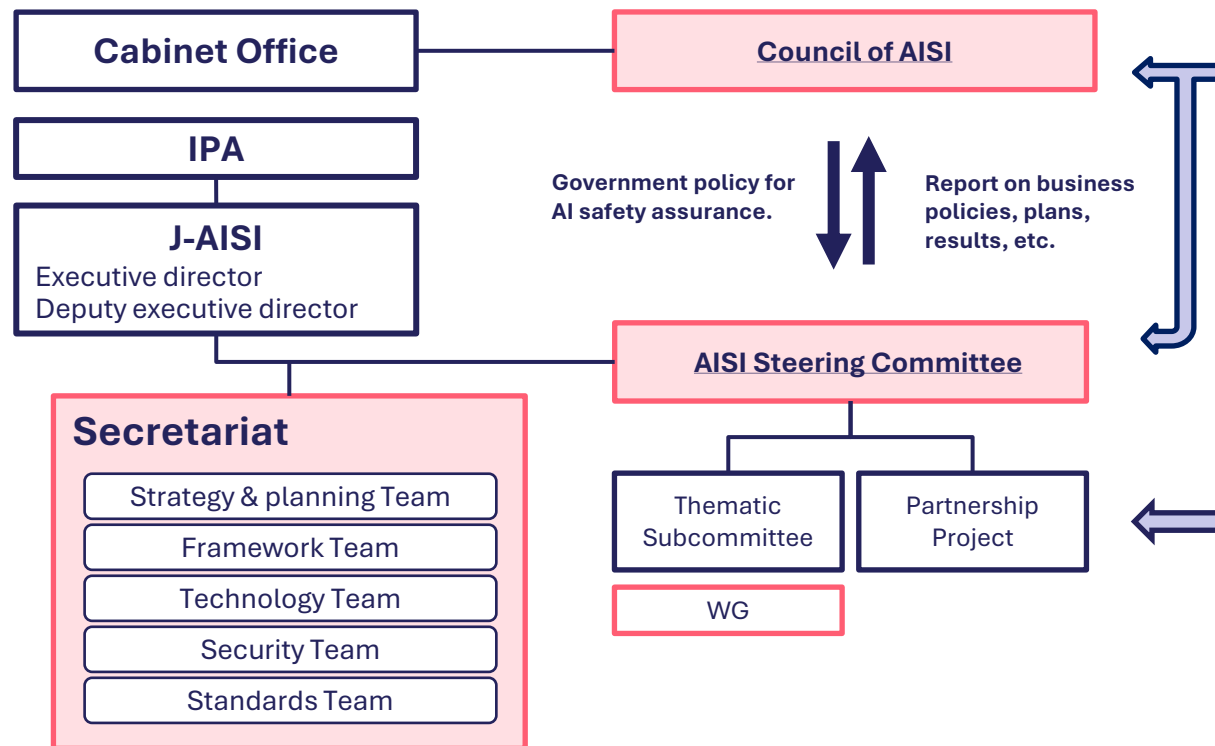
# Related Government organization and agencies

**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



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

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 for FY2024

		International	J-AISI	Government
		EVENT	DELIVERABLE	
2024	Apr		<ul style="list-style-type: none"><li>JP-U.S. Crosswalk1(4/30)</li></ul>	<ul style="list-style-type: none"><li>AI Guidelines for Business was published(4/19)</li></ul>
	May	AI Safety Summit, Korea		
	Jun	G7 Summit, Italy		<ul style="list-style-type: none"><li>Integrated Innovation Strategy 2024 was published(6/4)</li></ul>
	Jul		<ul style="list-style-type: none"><li>Japanese Translation of U.S. AI RMF(7/4)</li></ul>	
	Aug		<ul style="list-style-type: none"><li>Guide to Evaluation Perspectives(9/18)</li></ul>	
	Sep		<ul style="list-style-type: none"><li>JP-U.S. Crosswalk2(9/18)</li><li>Guide to Red Teaming Methodology*(9/25)</li></ul>	
	Oct			
	Nov	International Network of AISIs Convening, USA		
	Dec			
	Jan		<ul style="list-style-type: none"><li>Published Activity Map on AI Safety(2/7)</li><li>Published Data Quality Management Guidebook(Draft) (2/7)</li></ul>	
2025	Feb	AI Action Summit, France	<ul style="list-style-type: none"><li>Published National Status Report on AI Safety in Japan 2024(2/7)</li></ul>	<ul style="list-style-type: none"><li>Updated on AI Guidelines for Business (3/28)</li></ul>
	Mar			

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

## **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

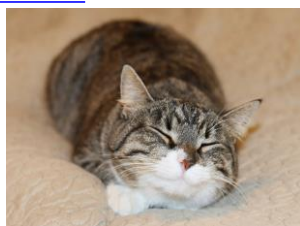
## Traditional AI

(Machine Learning)

### Classification

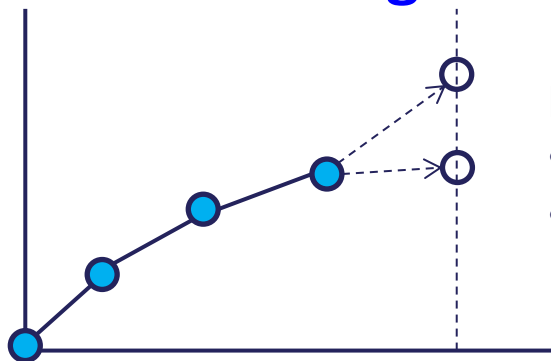


Dog or Cat → **Dog**



Dog or Cat → **Cat**

### Regression



#### Example

- Stock forecast
- Weather forecast etc

**The output format is simple**

## Generative AI

### 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

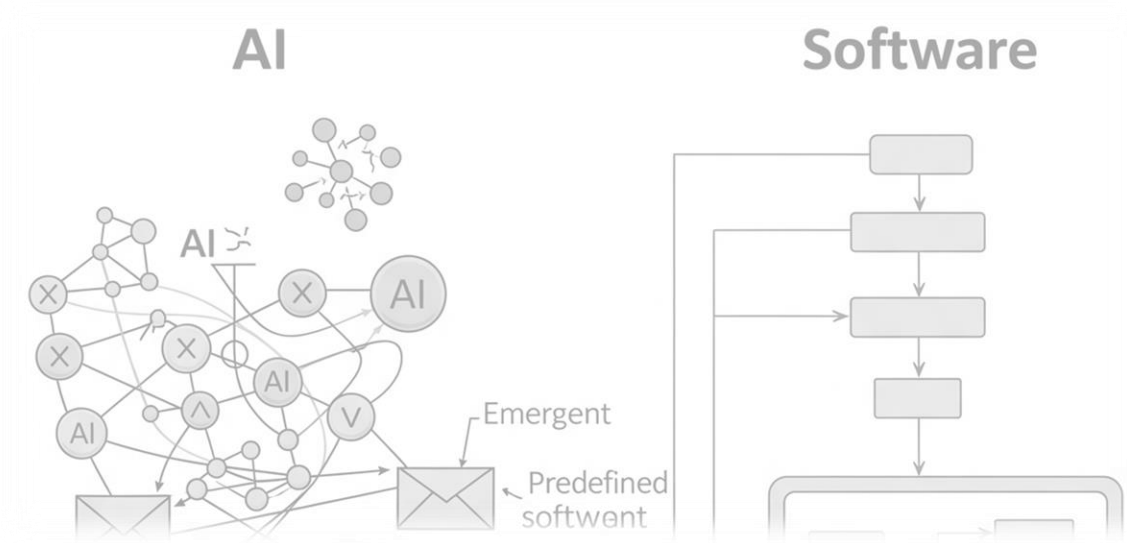
football on a  
grass



**Generate complex content**

# AI cannot be treated like software

- ◆ With software, the logic is defined, it's 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 detailed process of decision-making, which makes it unmanageable and causes problems.
  - Fine tuning of AI model and software patches are different.



## 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.(AI generated images must be used with care)

Fake Image: Mount Fuji Erupting



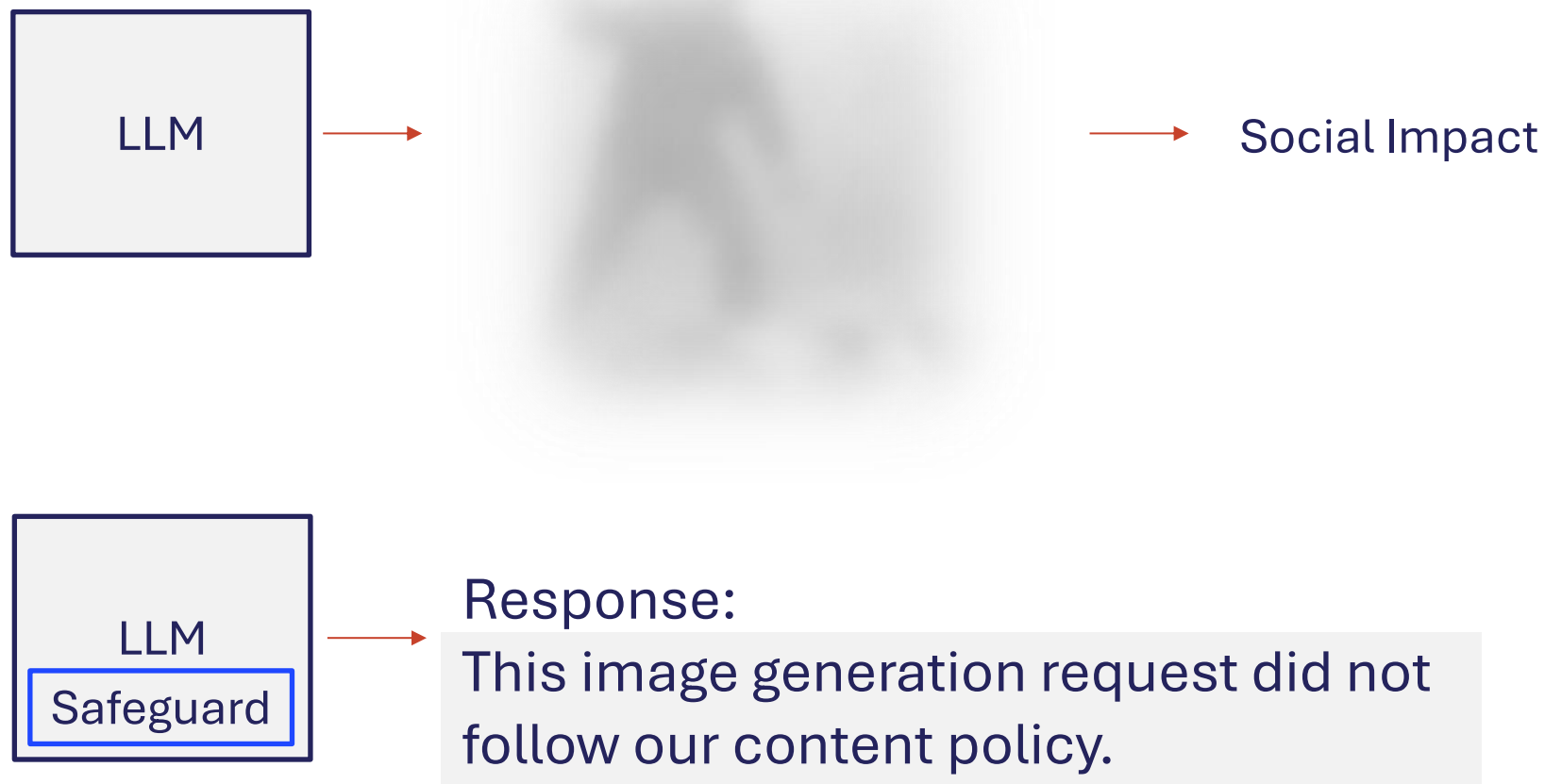
Fake Image: Tokyo Station Flooded



**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.



Introduction of **safeguards**  
as **AI safety** measures

## **3.2 How to achieve safe and secure AI**

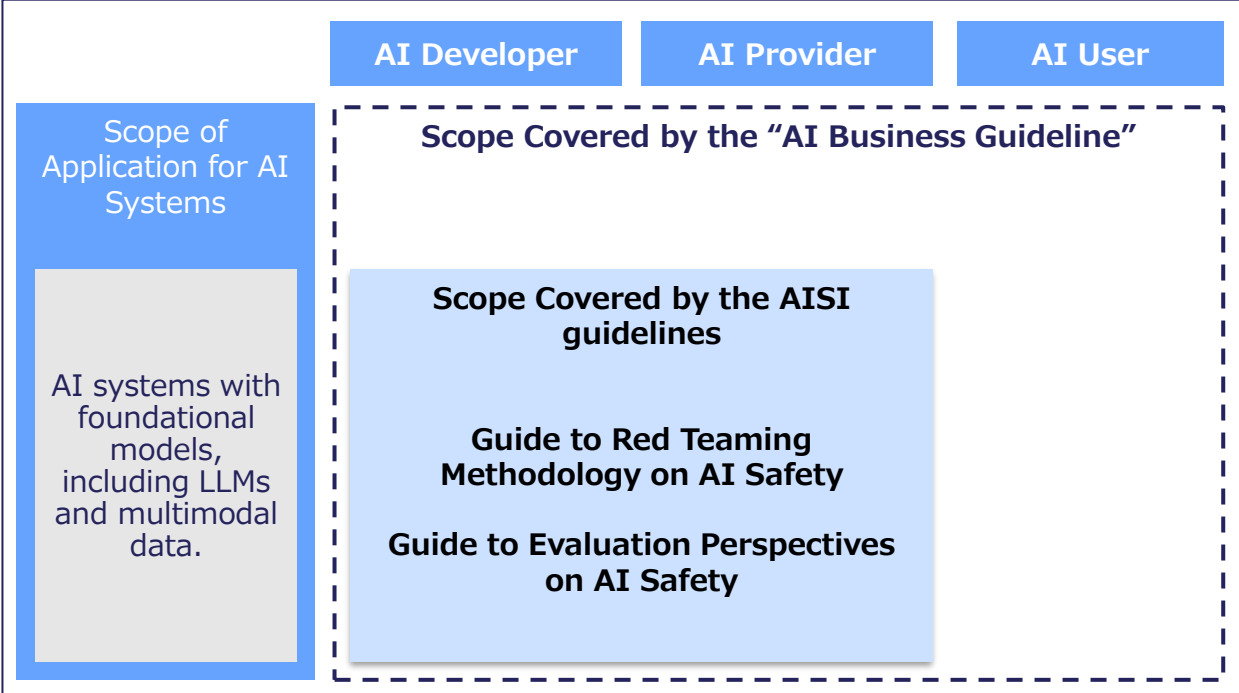
# Introduction to Key Guidelines for Realizing AI Safety

- 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.

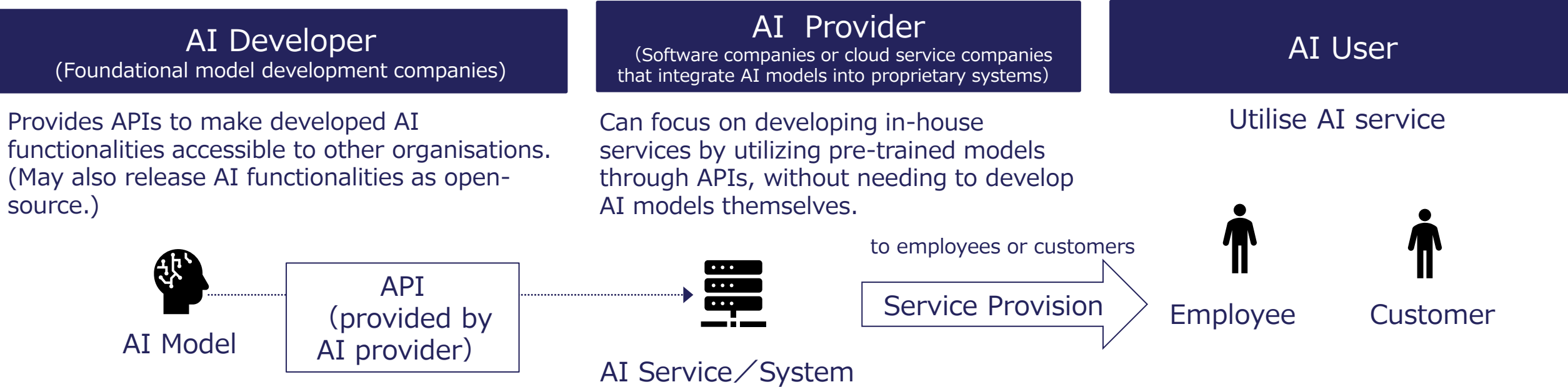
Key guidelines Related to AI Safety

Publisher	Title	Overview
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.
AISI	Guide to Red Teaming Methodology on AI Safety	A basic framework for considerations when conducting AI safety evaluations.
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.

Intended Audience of Each Guideline



## Flow from AI Development to Utilization (Reference)



AI Developer	AI capabilities (API)	Key Features
OpenAI	ChatGPT API	Generates text in conversational format.
Google	Google Cloud Natural Language API	Analyzes sentiment in text.
Stability AI	Stable Diffusion API	Generates images from text.

Examples of AI-Driven Services/Systems:
• Customer support chatbots.
• Tools for analyzing and categorizing review site posts.
• Character design tools.

End Users (Use Cases):
Customers: Asking questions through chatbots about products.
Employees: Using review analysis results for marketing.
Employees: Exploring character designs for advertising.

## 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.



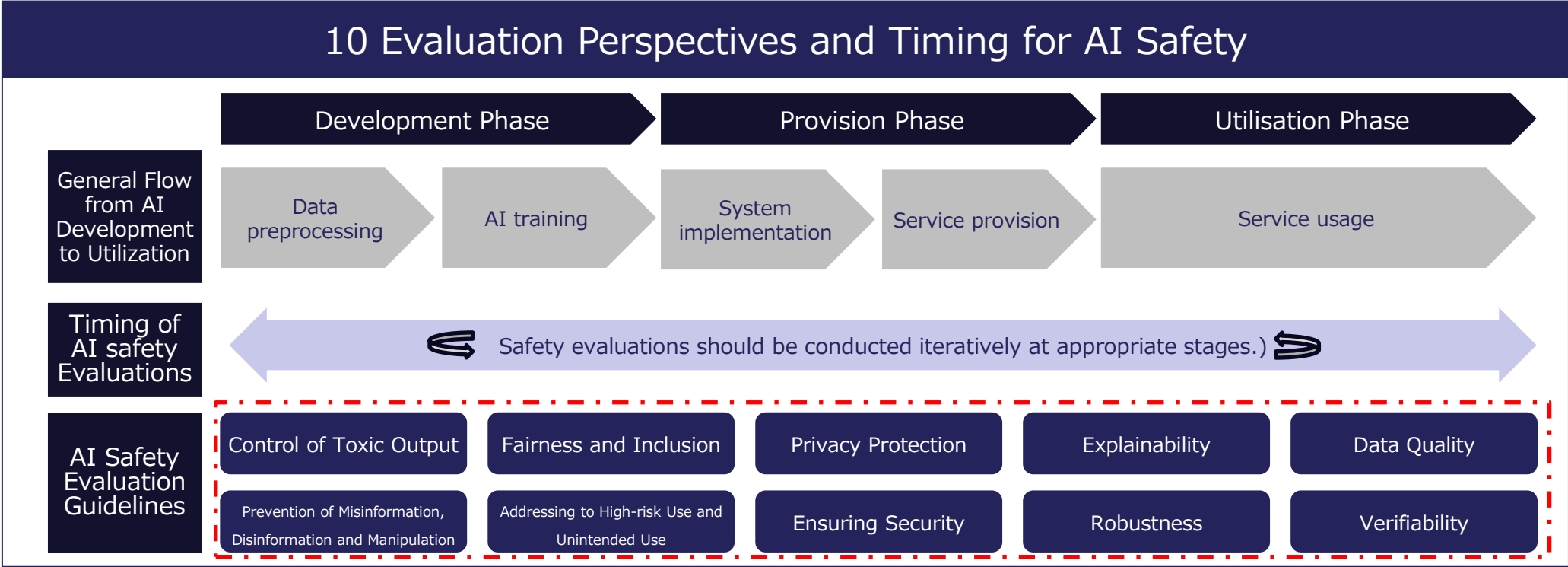
# What is AI Safety?



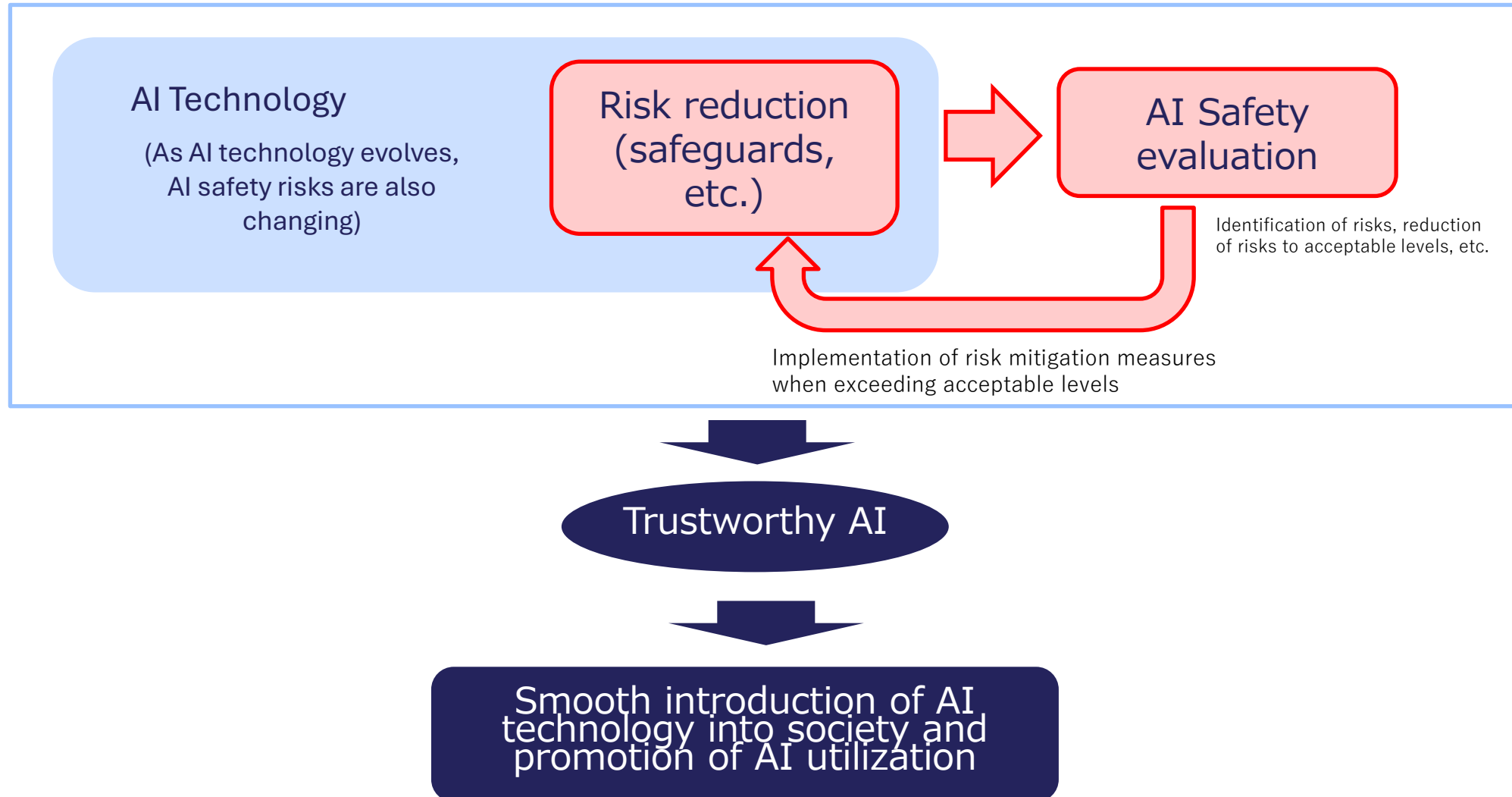
- ◆ AI Safety describes:
  - Based on a human-centric approach, it refers to a state in which
    - safety and fairness are maintained to reduce social risks\* associated with the use of AI,
    - privacy is protected to prevent inappropriate use of personal information,
    - security is ensured to 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**.

## 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.

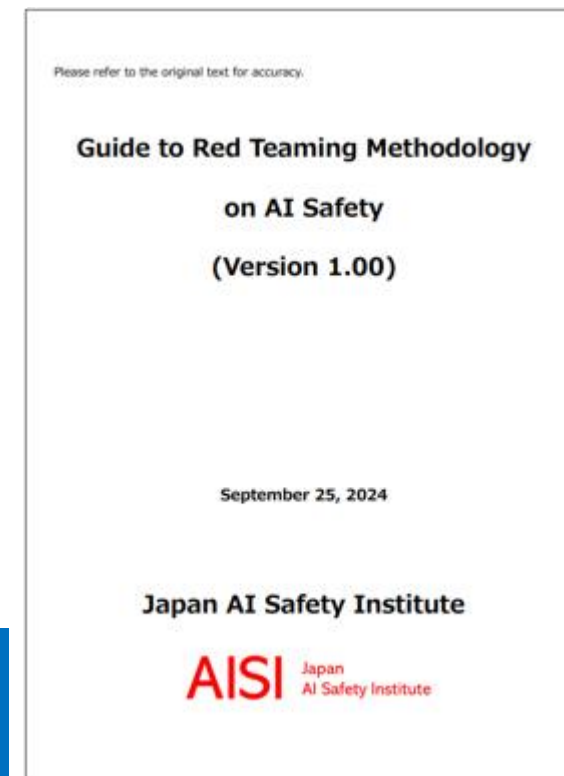


# Mitigation by AI safety evaluation



- ◆ 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: LLM System



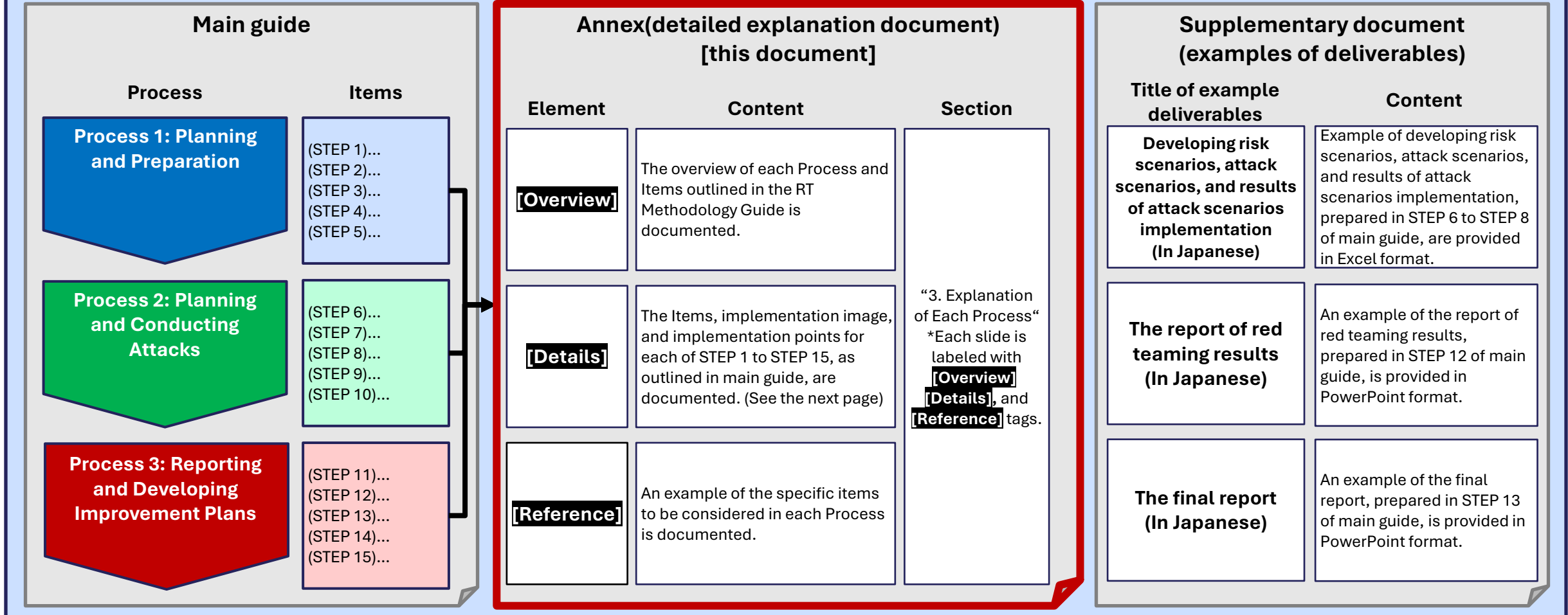
## Scope of Application for AI Systems (LLM System)

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

## 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.

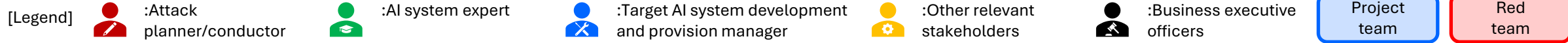
### Guide to Red Teaming Methodology on AI Safety



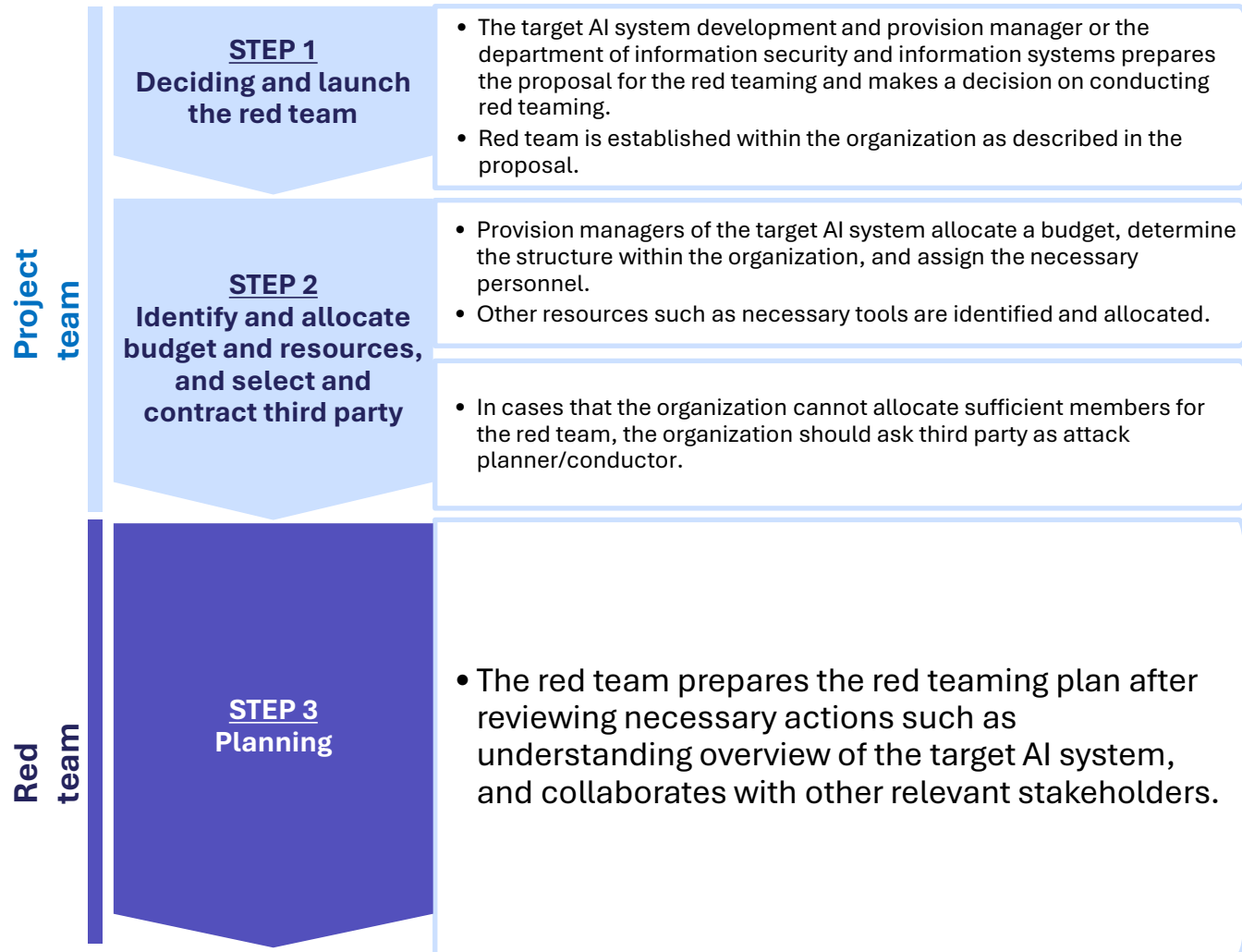
### 3. Explanation of Each Process Process 1



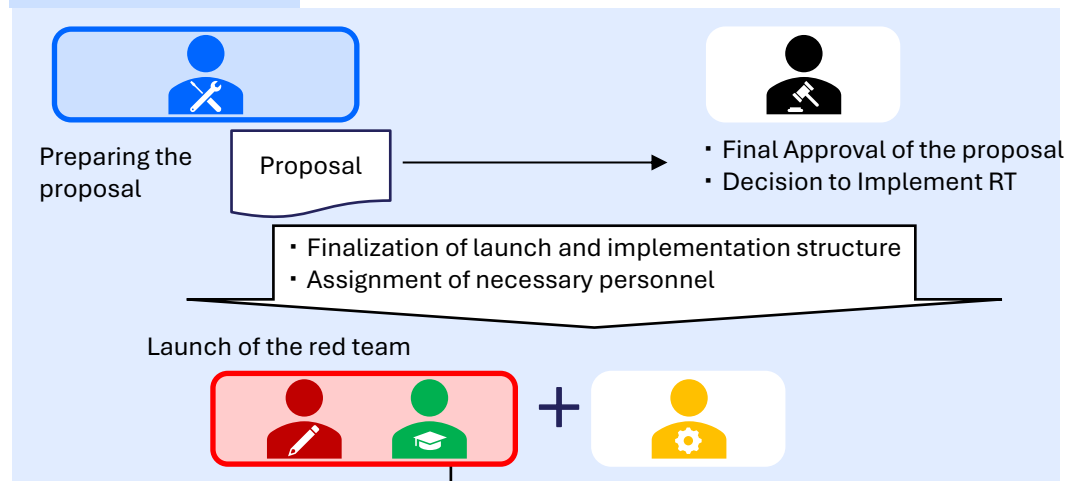
#### [Overview](STEP 1) Launch the team~(STEP 3) Planning



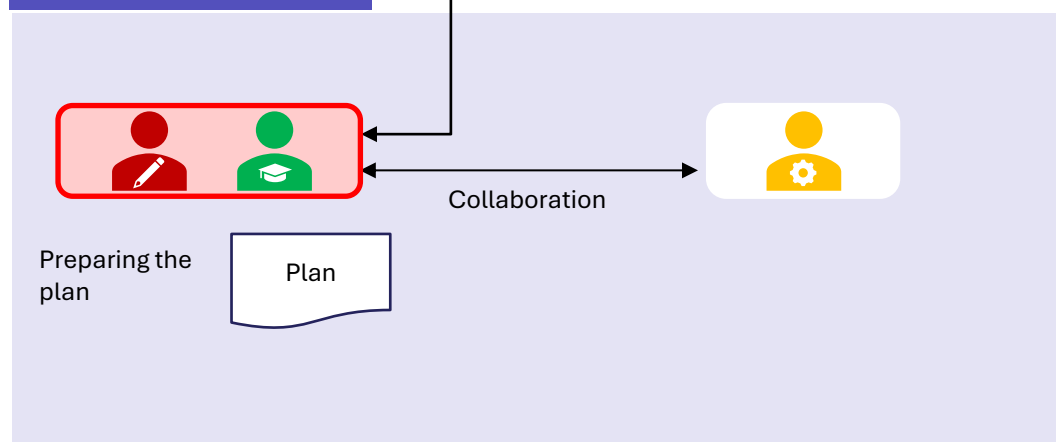
## Process 1: Planning and Preparation



### Project team



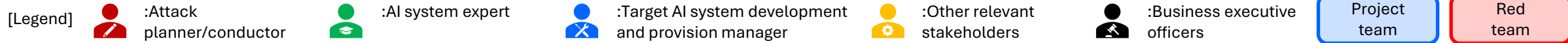
### Red team



### 3. Explanation of Each Process Process 1



#### [Overview](STEP 4) Preparing the environment, (STEP 5) Confirming escalation flow



## Process 1: Planning and Preparation

Red team

### STEP 4 Preparing the environment for red teaming

- Prior to conducting red teaming, the content, scope of impact, schedule, and other relevant details should be communicated to stakeholders.
- As needed, stakeholders are informed in advance and requested to temporarily disable monitoring settings, exclude themselves from monitoring, or ignore alerts.

### STEP 5 Confirming escalation flow

- The red team confirms escalation flow in case of unexpected behavior or failure/trouble due to red teaming conducting.

#### Red team



Collaboration



- Preparing the environment for red teaming
- Advance notification to stakeholders regarding RT content, scope of impact, and schedule

#### Red team



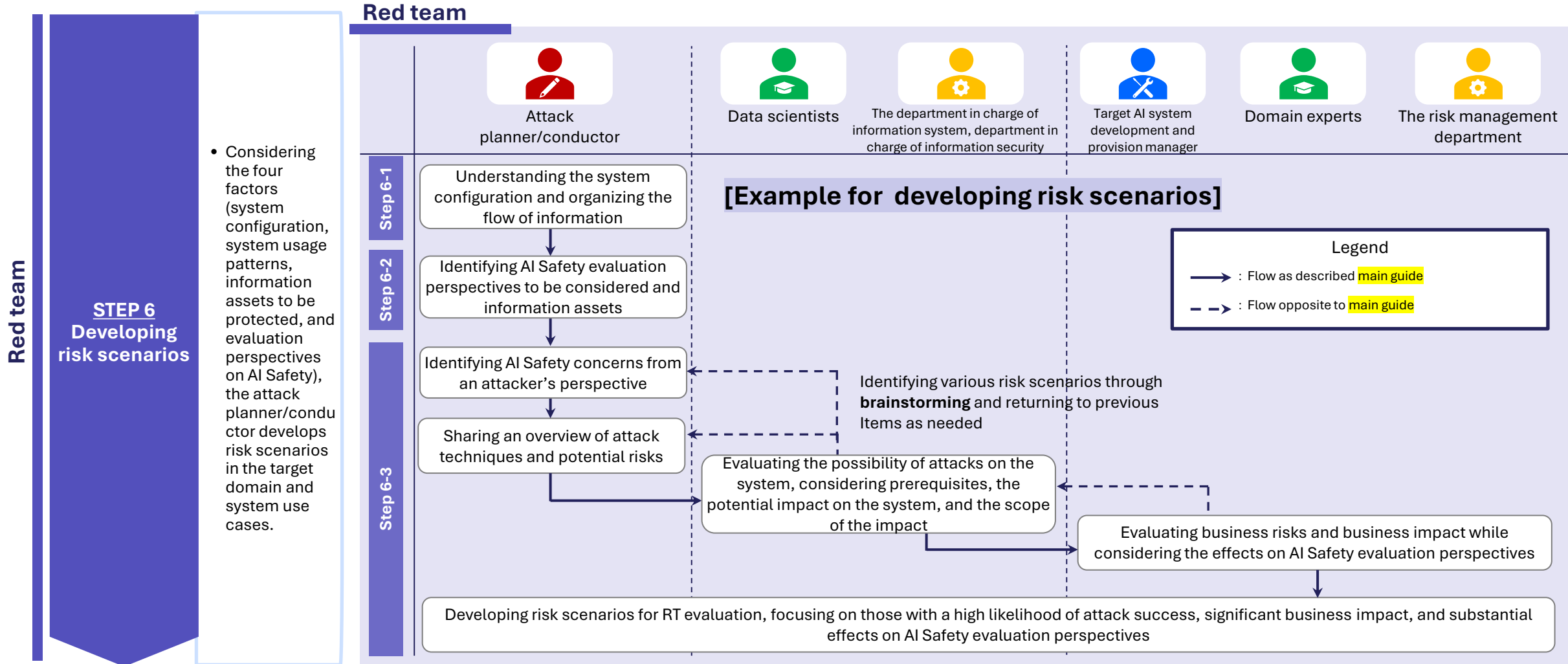
Collaboration



Confirming the escalation flow

#### [Overview](STEP 6) Developing risk scenarios

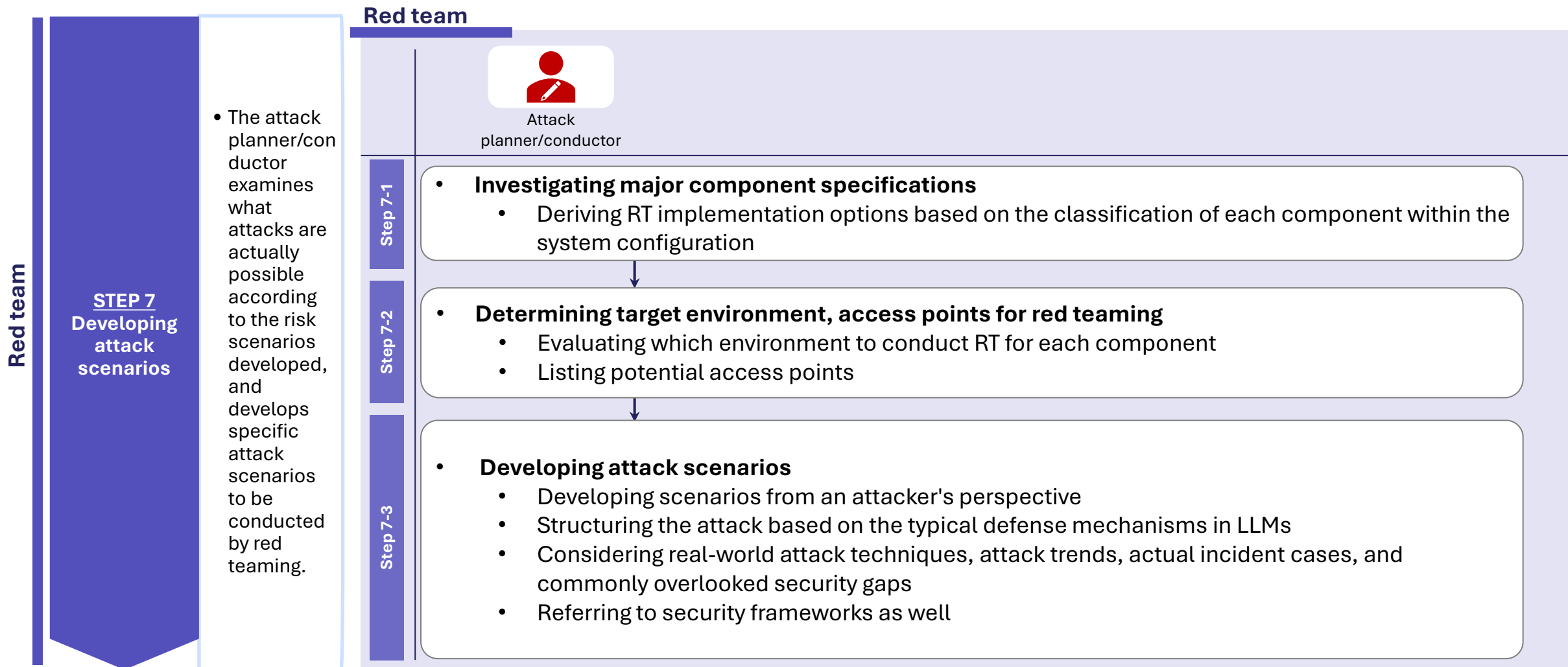
## Process 2: Planning and Conducting Attacks





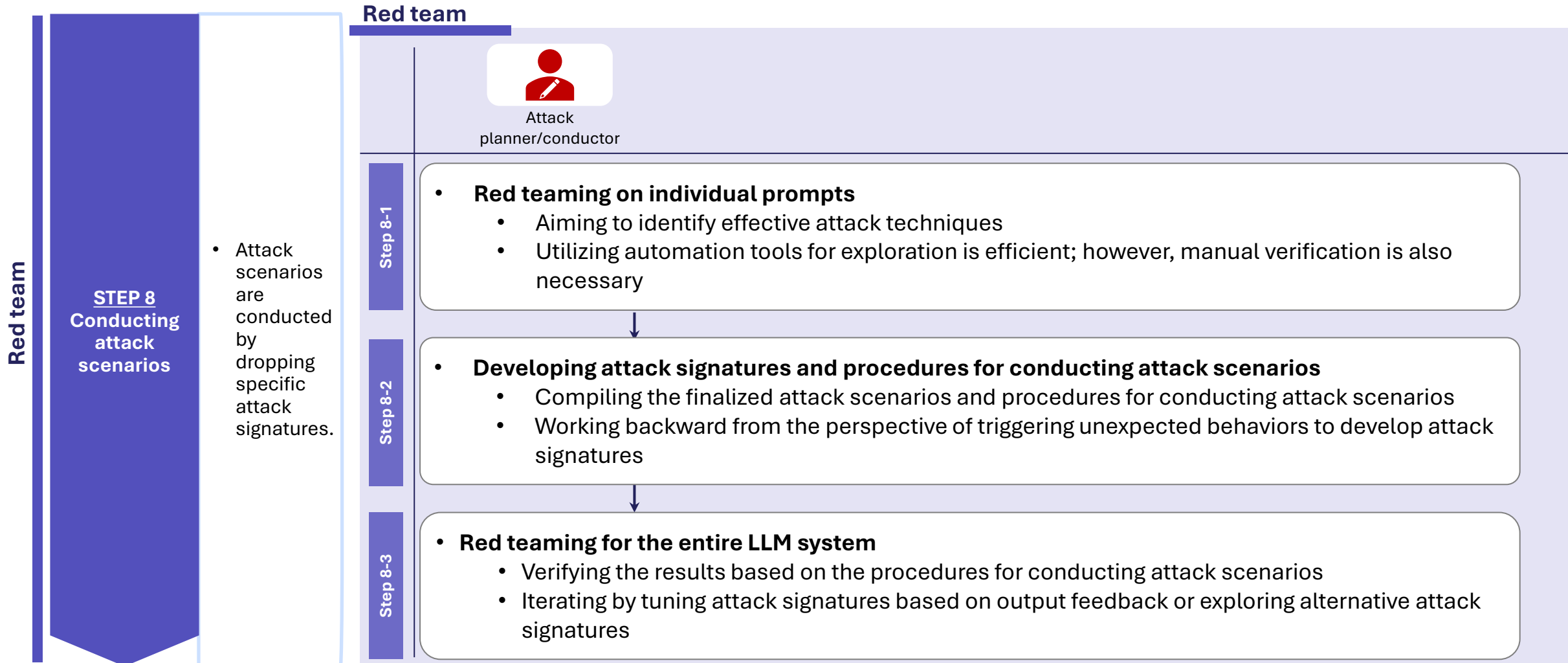
#### [Overview](STEP 7) Developing attack scenarios

## Process 2: Planning and Conducting Attacks



#### [Overview](STEP 8) Conducting attack scenarios

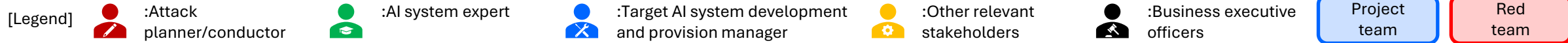
## Process 2: Planning and Conducting Attacks



### 3. Explanation of Each Process Process 2



#### [Overview](STEP 9) Record keeping, (STEP 10) After conducting attack scenarios



## Process 2: Planning and Conducting Attacks

Red team

### STEP 9

#### Record keeping during red teaming

- Records of red teaming in progress are kept in order to maintain a trail of the details of the red teaming conducted.

### STEP 10

#### After conducting attack scenarios

- The attack planner/conductor notifies the stakeholders, such as the development and provision managers of the target AI system and the department of information systems and information security, that red teaming attacks are finished.
- The temporary account for red teaming is deleted, and the settings are restored if any defensive measures that temporarily alter or relax the system settings have been implemented.

### Red team

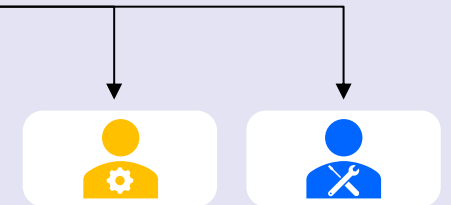


- Obtaining record during red teaming, documenting it in the report, and sharing it with relevant stakeholders
- For manual RT, capturing all attack signatures
- If an attack is successful, taking screenshots as evidence

### Red team



- Requesting the following actions from relevant stakeholders:
  - Notifying them of RT completion
  - Deleting temporary accounts created for RT
  - Restoring any configuration change



### 3. Explanation of Each Process Process 3

Process 1

Process 2

Process 3

11

12

13

14

15

[Overview](STEP 11) Analyzing the results, (STEP 12) Preparing the report of red teaming and review

[Legend]



:Attack  
planner/conductor



:AI system expert



:Target AI system development  
and provision manager



:Other relevant  
stakeholders



:Business executive  
officers

Project  
team

Red  
team

## Process 3: Reporting and Developing Improvement Plans

Red team

### STEP 11 Analyzing the red teaming results

- The attack planner/conductor analyzes the results obtained from red teaming.
- If necessary, additional confirmation of the information to be analyzed is made with relevant departments, such as the development and provision manager of the target AI system, the department in charge of information systems, and the department in charge of information security.

### STEP 12 Preparing the report of red teaming results and implementing stakeholder review

- Based on the vulnerabilities discovered during the red teaming exercise, the attack planner/conductor prepares logs and trails to prepare the overview of RT.
- The attack planner/conductor prepares the report of red teaming results and review it for factual errors, as necessary, with provision manager of the target AI system and with other relevant stakeholders.

### Red team



- Analyzing the implementation results
- Conducting additional verification and discussions as needed



### Red team



- Collecting and organizing logs and evidence
- Preparing the overview of RT
- Preparing the report of red teaming results and implementing stakeholder review

report of  
results

Stakeholder reviewing and reporting



### 3. Explanation of Each Process Process 3

Process 1

Process 2

Process 3

11

12

13

14

15

#### [Overview](STEP 13) Preparing and reporting the final results ~ (STEP 15) Follow-up

[Legend]



:Attack  
planner/conductor



:AI system expert



:Target AI system development  
and provision manager



:Other relevant  
stakeholders



:Business executive  
officers

Project  
team

Red  
team

## Process 3: Reporting and Developing Improvement Plans

Project team

### STEP 13 Preparing and reporting the final results

- The development and provision managers of the target AI system should prepare a final report of the red teaming, based on the report of red teaming results reported by the attack planner/conductor.
- If necessary, present the final report to the management team.

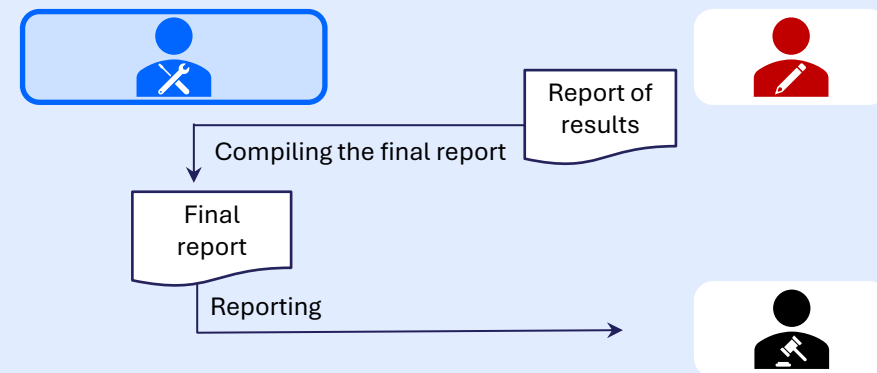
### STEP 14 Developing and implementing improvement plans

- The provision manager of the target AI system prepares improvement plans, specifying improvement measures to address business risks and other factors.
- When preparing improvement plans and measures, the project team should determine priorities based on the level of urgency and risk.

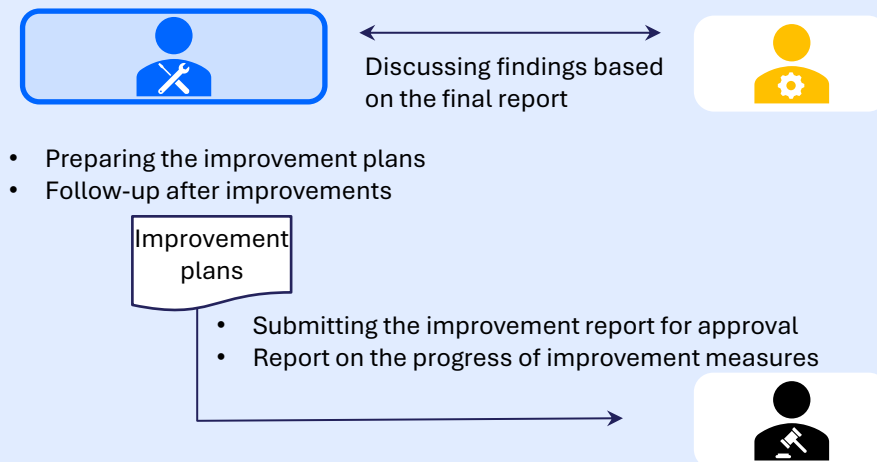
### STEP 15 Follow-up after improvement

- The progress of measures implemented based on the improvement plans should be checked at management meetings as appropriate.
- After implementing improvements measures, it is advisable to check the status of measures, review documents, or conduct red teaming again if necessary, to confirm that the vulnerability has been properly addressed and the risk has been mitigated.

### Project team



### Project team



# Real Sample Practice

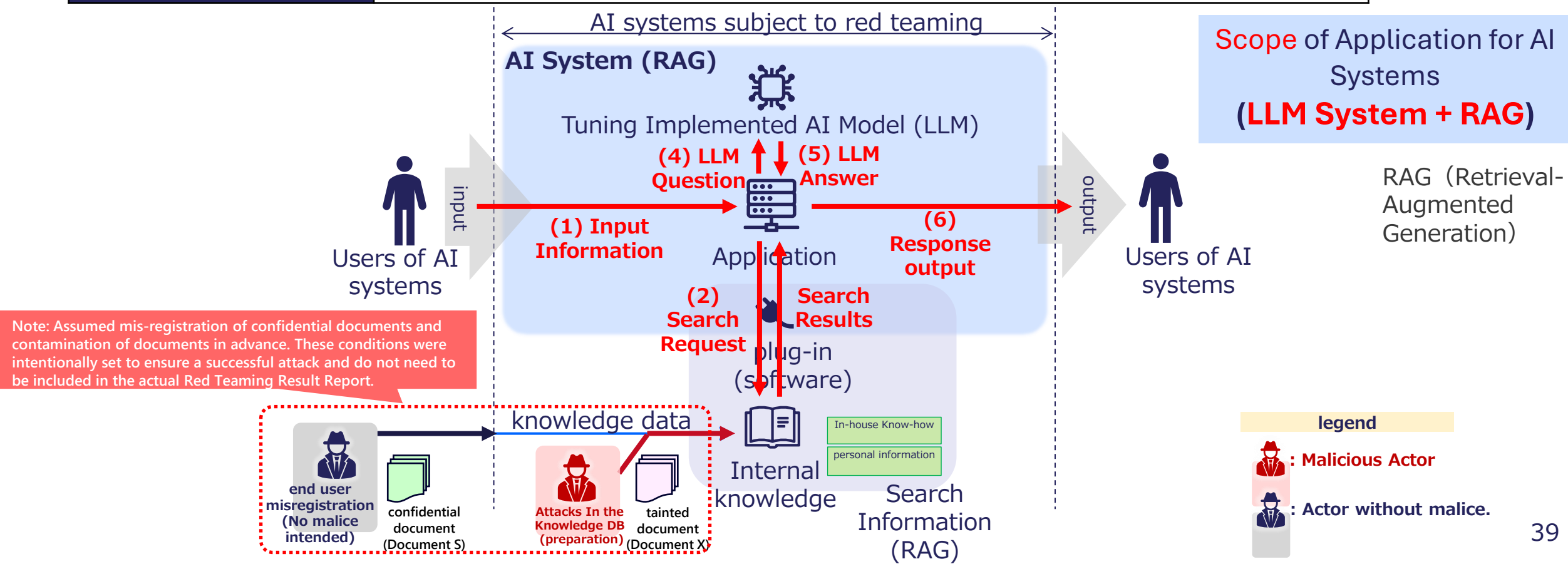
## System configuration diagram and information assets to be protected

### Outline of Target System

Based on the input information from the user of the AI system, the **relevant internal knowledge is retrieved**. The input information from user of the AI system and the result of the internal knowledge retrieval are then passed to the LLM, which generates output information based on the two pieces of information, and user of the AI system obtains the output.

### Information assets to be protected

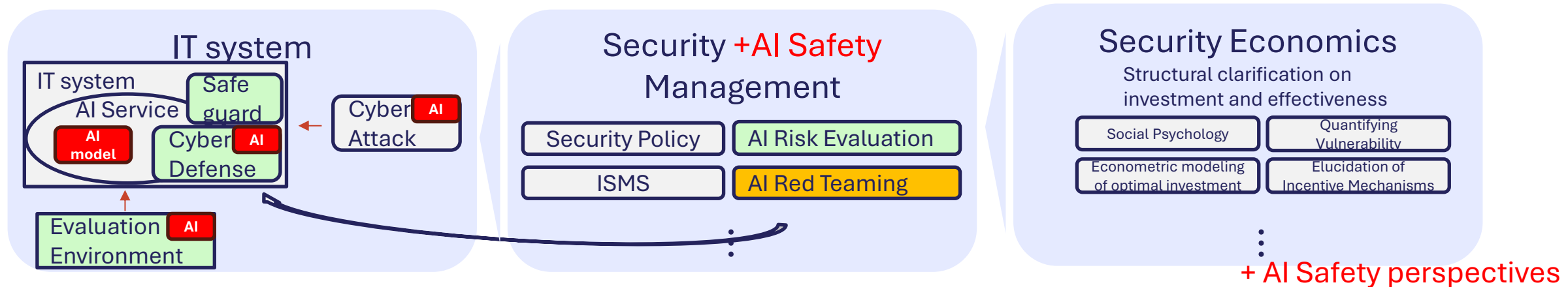
In-house know-how, personal information



## **5. Closing**



# AI Safety and Security Economics



# 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.

AISI Japan



Download guidelines here



Guide to Evaluation  
Perspectives on AI  
Safety



Guide to Red Teaming  
Methodology on AI Safety



AI Guidelines for  
Business



# AISI

Japan AI Safety Institute