

# **Guide to Red Teaming Methodology on AI Safety (Version 1.10)**

## **Summary**

**AI Safety Institute  
(March 31, 2025)**

# Table of Contents

<b>1. Background and Purpose</b>	<b>...P.3</b>
<b>2. Document Writing Policy</b>	<b>...P.4</b>
<b>3. Structure of This Document</b>	<b>...P.5</b>
<b>4. Scope of Red Teaming</b>	<b>...P.6</b>
<b>5. Overview of Red Teaming</b>	<b>...P.7</b>
<b>6. Process of Red Teaming</b>	<b>...P.9</b>
<b>Planning and Preparation</b>	
<b>Planning and Conducting Attacks</b>	
<b>Reporting and Developing Improvement Plans</b>	
<b>Revision of Guide to Red Teaming Methodology on AI Safety</b>	<b>...P.13</b>

## 1. Background and Purpose

**With the rapid proliferation of AI systems, AI Safety is becoming increasingly important. This document aims on presenting basic considerations for red teaming methodologies.**

### Background

- While the development, provision, and use of AI systems are expected to promote innovation and solve social problems, concerns have arisen due to misuse and abuse of AI systems, inaccurate output, etc.
- There is a growing interest in so-called AI Safety both in Japan and abroad, and **as part of the AI Safety evaluation, the red teaming method, in particular, is being studied in many countries.**

### Purpose

- The "Guide to Red Teaming Methodology on AI Safety" (hereafter referred to as "this document") provides **basic considerations for those involved in the development and provision of AI systems regarding red teaming methodologies to evaluate the risk countermeasures applied to the target AI system from an attacker's perspective.**
- This document presents items that are considered important for conducting red teaming at this stage, taking into consideration domestic and international studies, prior cases, and international alignment.

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."

\*Social risks include physical, psychological, and economic risks.

Source: <https://www.gov.uk/government/publications/ai-safety-institute-overview/introducing-the-ai-safety-institute>

## 2. Document Writing Policy

In addition to the AI Guidelines for Business, this document has been prepared based on a survey of domestic and international publications and tools.

### Domestic and International Publications

#### Machine Learning Quality Management Guidelines, 4th Edition (National Institute of Advanced Industrial Science and Technology)

The guidelines to classify and organize the quality requirements for AI systems utilizing machine learning in a top-down manner. They provide a structure that enables stakeholders involved in the target system to establish a framework for the objective evaluation of the system's quality.

#### LLM AI Cyber Security and Governance Checklist (Open Worldwide Application Security Project (OWASP))

It describes the management of cybersecurity risks in the provision and use of AI systems within an organization, including a list of the top 10 critical vulnerabilities in LLM applications.

#### SP800-115, National Institute of Standards and Technology (NIST)

Comprehensive guidelines for information systems security testing and evaluation.

#### AI 800-1(Initial Public Draft) (National Institute of Standards and Technology (NIST))

Draft guidelines for risk management of dual-use foundation models.

### AI Guidelines for Business(Japan)

Guidelines developed by integrating and updating existing relevant guidelines in Japan in order to respond to rapid technological changes in recent years.

## Guide to Red Teaming Methodology on AI Safety

### Tools (Organization)

#### Anthropic

Anthropic has publicly announced that it conducts red teaming from various perspectives and publishes the data sets for red teaming.

#### Microsoft

Microsoft is conducting red teaming for its services and has published a guide on red teaming.

#### NVIDIA

Nvidia has conducted red teaming by cross-disciplinary teams within its own company.

#### OpenAI

OpenAI is recruiting experts to red team their AI model, working to enhance the safety of their products.

#### Project Moonshot (AI Verify Foundation)

Project Moonshot is an open-source tool developed by the AI Verify Foundation in Singapore, with features to support red teaming exercise.

### 3. Structure of This Document

Items considered important for conducting red teaming on AI Safety are categorized by type. The table of contents is organized according to the categories to enhance readability.

- The contents of each section of this document are described based on the items organized from a 5W1H perspective.
- The primary target audience is assumed to be AI developers and AI providers. In particular, the target readers are “development and provision managers” and “business executive officers” who are involved in the planning and conducting red teaming.

Type	Examples of items to be described
<b>What</b> (What is red teaming?)	<ul style="list-style-type: none"> <li>➤ Definition and scope of “red teaming”</li> <li>➤ AI systems covered in this publication</li> </ul>
<b>Why</b> (Why red teaming?)	<ul style="list-style-type: none"> <li>➤ Purpose of red teaming</li> <li>➤ Importance and expected effects of red teaming</li> </ul>
<b>Who</b> (Who will conduct red teaming?)	<ul style="list-style-type: none"> <li>➤ What roles are the red teaming conductors?</li> </ul>
<b>When</b> (When to conduct red teaming?)	<ul style="list-style-type: none"> <li>➤ Timing of red teaming</li> </ul>
<b>Where</b> (where to conduct red teaming?)	<ul style="list-style-type: none"> <li>➤ Whether it will be performed by your own organization or by a third party</li> </ul>
<b>How</b> (How to conduct red teaming?)	<ul style="list-style-type: none"> <li>➤ How to plan red teaming and what to prepare for it</li> <li>➤ What threats to assume in red teaming</li> </ul>

Main guide to Red Teaming Methodology on AI Safety [Table of Contents]	
1	Introduction
2	About Red Teaming
3	Typical Attack Methods on LLM systems
4	Red Teaming Structure and Roles
5	Timing of Red Teaming and its Process
6	Planning and Preparation
7	Planning and Conducting Attacks
8	Reporting and Developing Improvement Plans
A	Appendix



**Intended Audience**

AI Developers and AI Providers

Development and Provision Managers



Business Executives Officers



\*Readers who are involved in the planning and conducting of red teaming, among those listed on the left.

In the preparation of Version 1.10, Annex (detailed explanation document) and Supplementary document (examples of deliverables) were prepared in addition to main guide. For more details, please refer to page 15 of this document.

## 4. Scope of Red Teaming

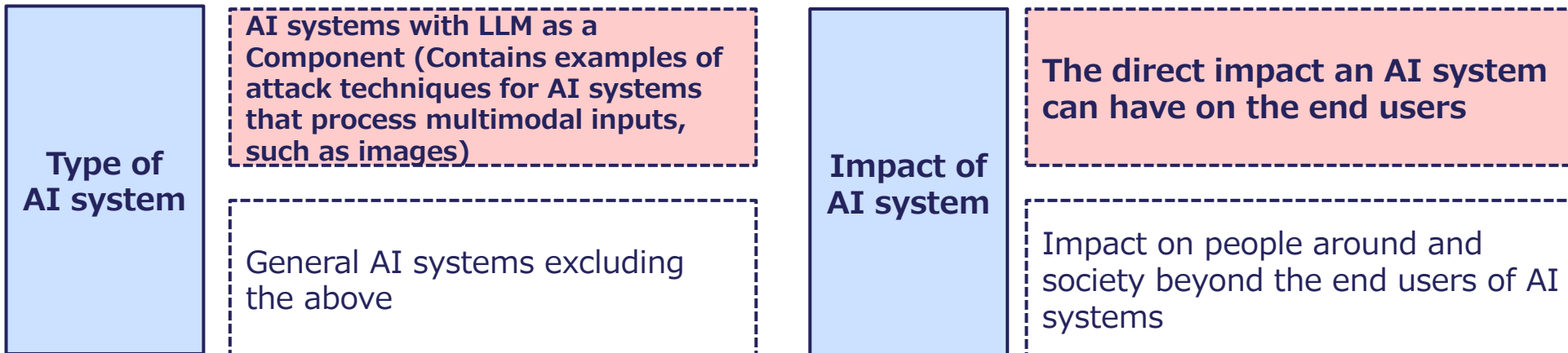
Red teaming in this document describes as an evaluation method to check the effectiveness of response structure and countermeasures for AI Safety in terms of how attackers attack AI systems. The red teaming described here is specifically focused on AI systems(LLM systems) that incorporate large language models.

### 🔍 Scope of red teaming in this document

- **Red teaming is “an evaluation method to check the effectiveness of response structure and countermeasures for AI Safety in terms of how attackers attack AI systems.”** In this document, red teaming with respect to AI Safety is simply referred to as "red teaming".
- The scope of the red teaming in this document is organized in terms of **(1) Type of AI system** and **(2) Impact of AI system**. (Red teaming is one of the AI Safety evaluation methods, and its scope is consistent with the one described in the “Guide to Evaluation Perspectives on AI Safety.”)

#### Scope of red teaming in this document

The boxes noted in red below indicate the scope of the red teaming in this document.



## 5. Overview of Red Teaming

Red teaming is "an evaluation method to check the effectiveness of response structure and countermeasures for AI Safety in terms of how attackers attack AI systems," and is one of the methods of AI Safety evaluation.

- The purpose of red teaming is to maintain or enhance AI safety by identifying vulnerabilities such as weaknesses and insufficient countermeasures in the target AI system from an attacker's perspective, then mitigating them through system hardening.

### Types of Red Teaming

- Red teaming can be categorized as follows.

#### Category of red teaming tests based on prior knowledge of the attack planner/conductor

- **Black-box Test**  
(The attack planner/conductor does not have any prior knowledge of the system, such as its internal structure.)
- **White-box Test**  
(The attack planner/conductor has sufficient knowledge of the system, such as its internal structure.)
- **Gray-Box Test**  
(The attack planner/conductor has partial knowledge of the system, such as its internal structure.)

#### Category of the environment in which red teaming is conducted

- **Production Environment**  
(Production environment where AI systems are actually put into practice)
- **Staging Environment**  
(Environment for testing and checking for defects in conditions similar to those of the actual production environment)
- **Development Environment**  
(Environment for developing AI systems)

#### Category of how attack signatures are attempted

- **Red Teaming with Automated Tools**
- **Manual Red Teaming**
- **Red Teaming with AI Agents**

### Typical attack methods on LLM systems

- Examples of typical attack methods against LLM systems. They should be considered in red teaming.
  - **Direct Prompt Injection**  
Attacker directly injects malicious prompts into the AI system
  - **Indirect Prompt Injection**  
Attacker indirectly injects malicious prompts into the AI system
  - **Prompt Leaking**  
Attacker extracts the designated system prompt
  - **Poisoning Attacks**  
Attacker infiltrates manipulated data or model into data or model during training
  - **Evasion Attacks**  
Malicious modification of inputs to the AI system to cause unintended behavior
  - **Model Extraction Attack**  
An attack to create a model with the same performance as the target system's model by analyzing its inputs and outputs
  - **Membership Inference Attacks**  
An attack that identifies whether certain data is included in the training data by analyzing system's inputs and outputs
  - **Model Inversion Attacks**  
An attack that recovers information contained in training data by analyzing inputs and outputs

## 5. Overview of Red Teaming

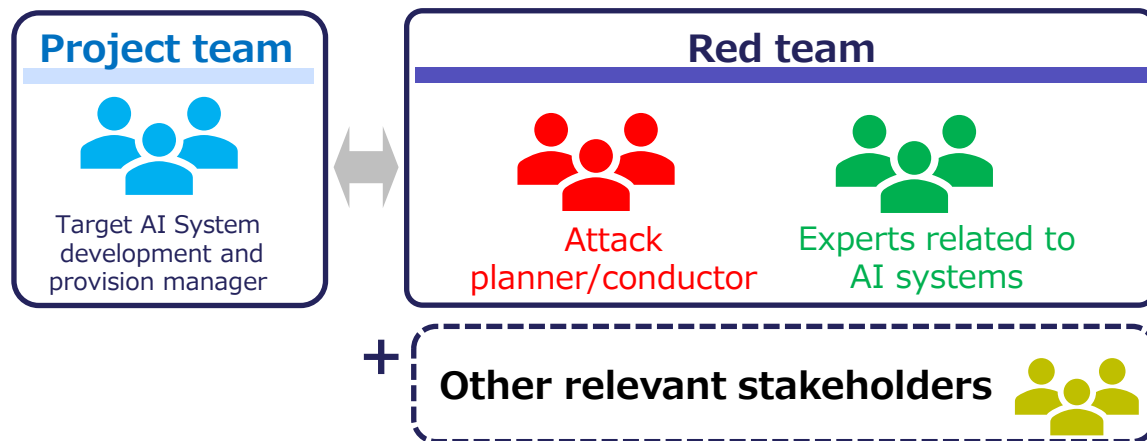
It is desirable to have diverse stakeholders participate in the conducting of red teaming. It is also desirable to conduct the red teaming not only before the release of the AI system, but also after the start of operation, as needed.

### Conducting Structure and Roles

- A red team\* should be established in coordination with the project team involved in the development and provision of the AI system subject to red teaming, and a leader or responsible person should be appointed.
- It is basically assumed that the red team includes the "Attack Planner/Conductor," and "Experts related to AI systems."
- Other relevant stakeholders within the organization may be involved as needed.

(\* ) Team in charge of checking the effectiveness of the response structure and countermeasures for AI Safety in terms of how attackers attack AI systems.

#### Red team structure



### Conducting Timing

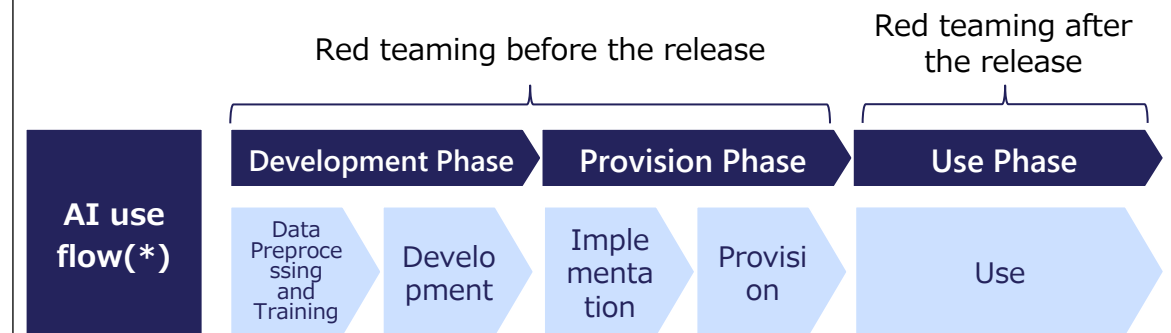
#### Red teaming before the release

- The first time red teaming is conducted, it should be done before the release of the target AI system.
- However, depending on the scale and complexity of the target AI system, it may be effective to conduct risk analysis in units such as system components and system layers, and divide and red teaming at appropriate timing.

#### Red teaming after the release

- Red teaming is not a one-time task; it should be conducted periodically as needed.

#### Timing of red teaming in AI use flow



(\*) Refer to the AI Guidelines for Business "Correlation between AI business actor and AI use flow".



## 6. Process of Red Teaming

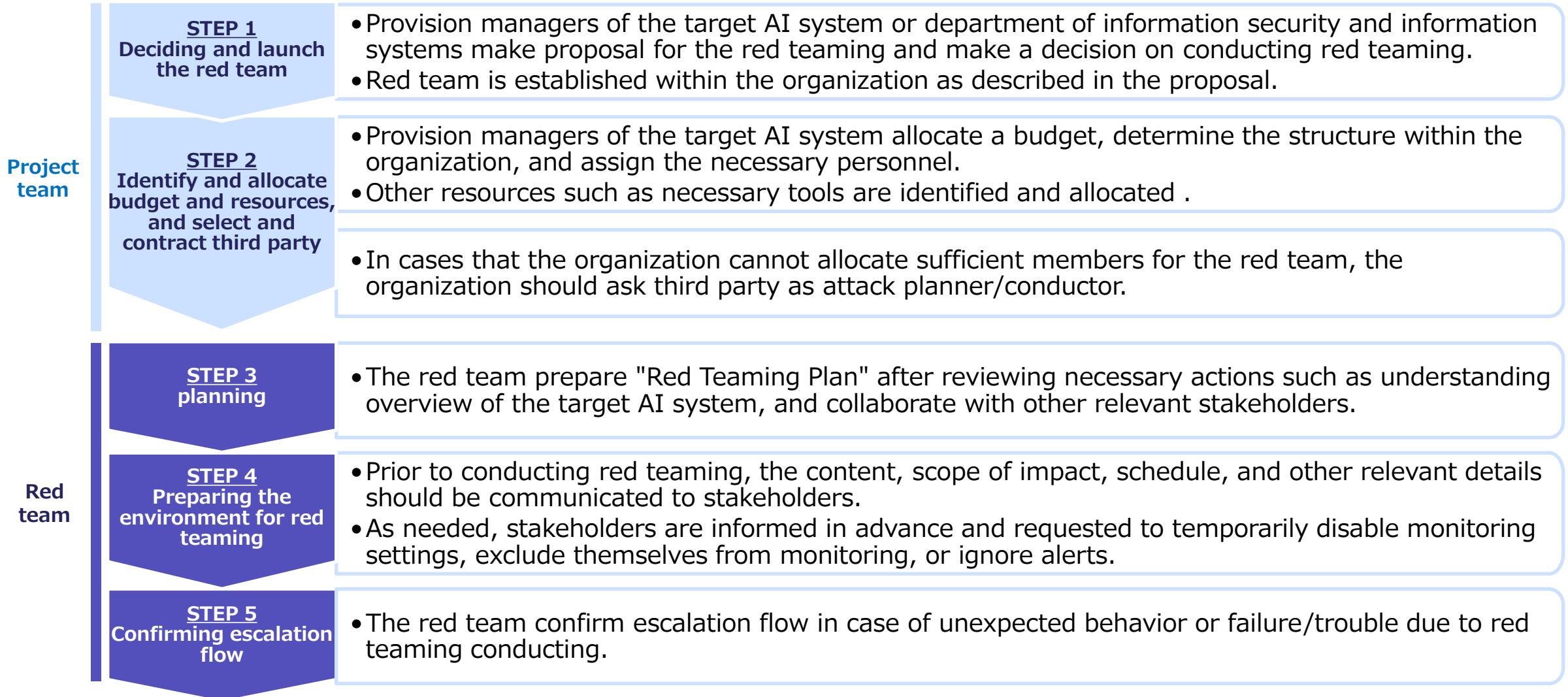
The red teaming Process consists of three parts: "Planning and Preparation," "Planning and Conducting Attacks," and "Reporting and Developing Improvement Plans."

Process	Items	Chapter in main guide
<b>Process 1: Planning and Preparation</b>	<ul style="list-style-type: none"><li>✓ Deciding and launch the red team</li><li>✓ Identify and allocate budget and resources, and select and contract third party</li><li>✓ Planning</li><li>✓ Preparing the environment for red teaming</li><li>✓ Confirming escalation flow</li></ul>	Chapter 6.
<b>Process 2: Planning and Conducting Attacks</b>	<ul style="list-style-type: none"><li>✓ Developing risk scenarios</li><li>✓ Developing attack scenarios</li><li>✓ Conducting attack scenarios</li><li>✓ Record Keeping during red teaming</li><li>✓ After conducting attack scenarios</li></ul>	Chapter 7.
<b>Process 3: Reporting and Developing Improvement Plans</b>	<ul style="list-style-type: none"><li>✓ Analyzing the red teaming results</li><li>✓ Preparing the report of red teaming results and implementing stakeholder review</li><li>✓ Preparing and reporting the final results</li><li>✓ Developing and implementing improvement plans</li><li>✓ Follow-up after improvement</li></ul>	Chapter 8.

## 6. Process of Red Teaming (Planning and Preparation)

In “Planning and Preparation,” the project team launches the red team, develop a red teaming plan, and take necessary actions such as allocating budgets.

### Process 1: Planning and Preparation



## 6. Process of Red Teaming (Planning and Conducting Attacks)

In “Planning and Conducting Attacks,” the red team develops risk and attack scenarios, conduct attack scenarios, keep the records, among other things.

### Process 2: Planning and Conducting Attacks

Red  
team

#### STEP 6 Developing risk scenarios

- Considering the four factors (system architecture, system usage patterns, information assets to be protected, and evaluation perspectives on AI Safety), the attack planner/conductor develops risk scenarios in the target domain and system use cases.

#### STEP 7 Developing attack scenarios

- The attack planner/conductor examines what attacks are actually possible according to the risk scenarios developed, and develop specific attack scenarios to be conducted by red teaming.

#### STEP 8 Conducting attack scenarios

- Attack scenarios are conducted by dropping specific attack signatures.

#### 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.

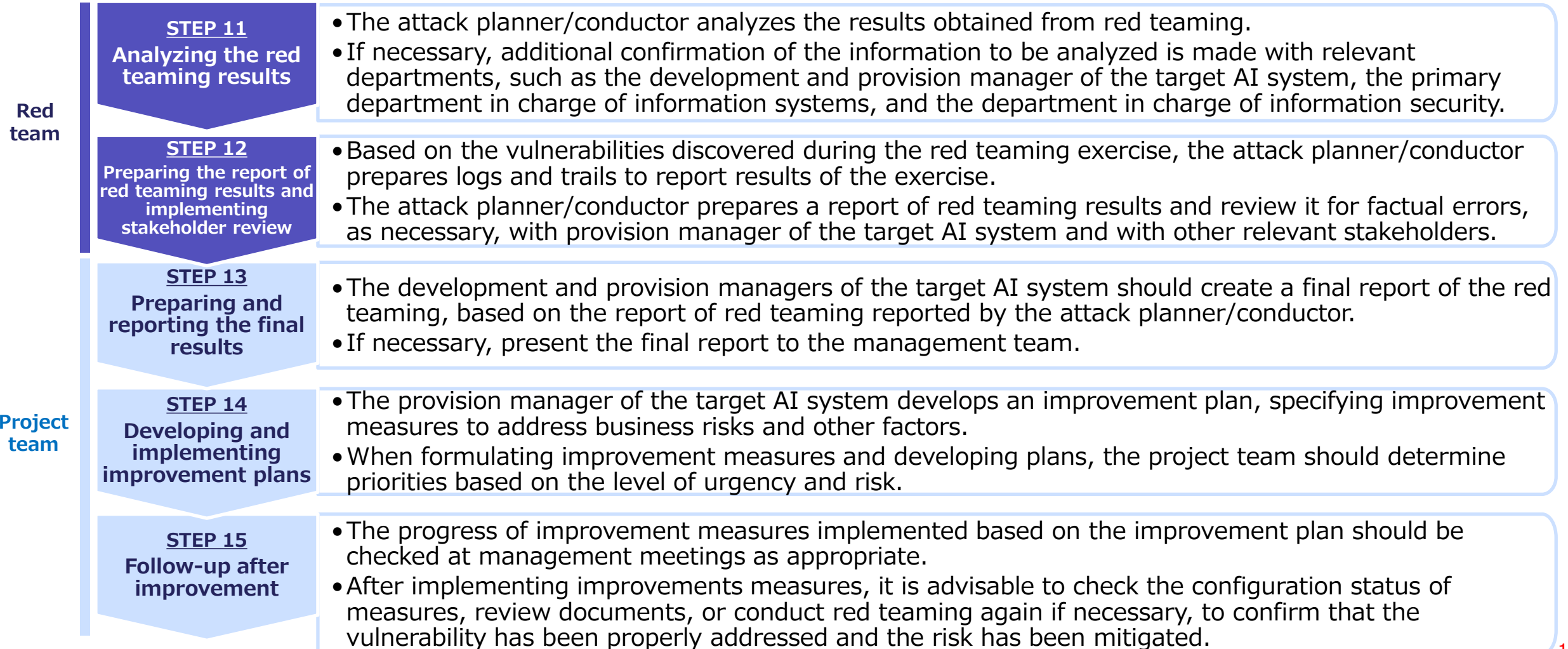
#### STEP 10 After conducting attack scenarios

- The attack planner/conductor notifies the stakeholders such as the development and provision managers of the target AI system and department of information systems and information security that attacks of red teaming 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.

## 6. Process of Red Teaming (Reporting and Developing Improvement Plans)

In "Reporting and Developing Improvement Plans," improvements are made to the items identified as a result of the red teaming. After the results of conducting are reviewed, an improvement plan is developed and implemented, and follow-up on the improvement measures is implemented.

### Process 3: Reporting and Developing Improvement Plans



# **Revision of Guide to Red Teaming Methodology on AI Safety**

**AI Safety Institute  
(March 31, 2025)**

**This document has been revised to accommodate a broader range of red teaming methodologies, making it a more practical guide for reference.**

### Background

- In September 2024, this document was prepared to provide those involved in the development and provision of AI systems with basic considerations regarding red teaming methodologies to evaluate the countermeasures to risk applied to a target AI system from an attacker's perspective.
- The document is systematically organized based on the "AI Guidelines for Business," as well as a review of literature and relevant stakeholders both domestically and internationally. However, given that Process 2 of red teaming (Planning and Conducting Attacks) requires a high level of expertise, there is a need to enhance the document to serve as a more practical guide.

### Purpose

- A study were conducted to expand the scope of evaluation of this document to include multimodal foundation models, and areas for revision were considered.
- Red teaming were conducted in accordance with this document, and revisions were considered to enhance it as a more practical guide. The focus was on use cases involving LLM system using RAG, likely to be adopted by Japanese companies.

\*1: AI models that are multimodal (capable of handling multiple modalities) and trained on extensive data, adaptable to a wide range of downstream tasks.

## Implementation details for revision of Guide to Red Teaming Methodology on AI Safety

Guide to Red Teaming Methodology on AI Safety as of September 2024 was used as main guide, with some of its contents added, and detailed explanation document as Annex and examples of deliverables as Supplementary document were newly prepared.

The revision points were examined if the evaluation scope is expanded to include multimodal foundation models.

[Revision contents]

Reference information on attack methods concerning multimodal information such as images was added to main guide under "Chapter 3. Typical Attack Methods on LLM Systems."

Based on the results of the implementation aligned with main guide, revisions were considered to make it a practical guide.

[Revision contents]

To help readers better understand and effectively practice each phase of red teaming, the guide as of September 2024 was established as main guide, with some content added, along with Annex that provides detailed explanation document and Supplementary document that provides examples of deliverables.

- ✓ Main guide: flowcharts and other visual aids were included to facilitate a clearer understanding of the red teaming flow for readers.
- ✓ Annex (detailed explanation document): This explains implementation points for conducting red teaming, enabling readers to apply the guide practically.
- ✓ Supplementary document (examples of deliverables) (In Japanese): Examples of deliverables were prepared to help readers visualize the expected deliverables in Japanese.

Add

Prepare new

### Guide to Red Teaming Methodology on AI Safety [this document]

#### Main guide

##### Main guide [Table of Contents]

1	Introduction
2	About Red Teaming
3	Typical Attack Methods on LLM systems
4	Red Teaming Structure and Roles
5	Timing of Red Teaming and its Process
6	Planning and Preparation
7	Planning and Conducting Attacks
8	Reporting and Developing Improvement Plans
A	Appendix

#### Annex (detailed explanation document)

##### Annex (Detailed Explanation Document) [Table of Contents]

1	Background and Purpose of the Detailed Explanation Document
2	Role of the Detailed Explanation Document
3	Explanation of Each Process

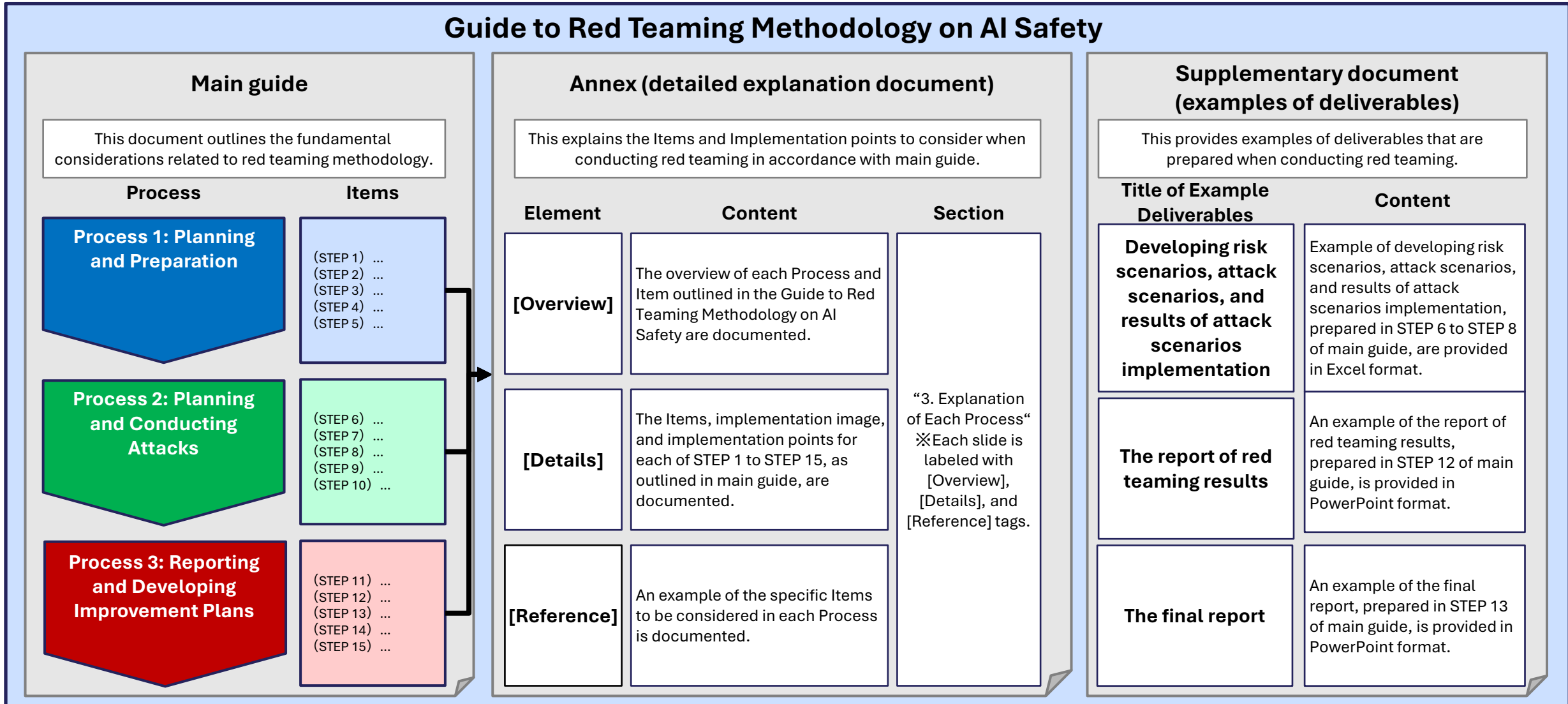
#### Supplementary document (examples of deliverables) (In Japanese)

##### Supplementary document (Examples of Deliverables)

1	Developing risk scenarios, attack scenarios, and results of attack scenarios implementation (Excel format)
2	The report of red teaming results (PowerPoint format)
3	The final report (PowerPoint format)

The Items in each Process of main guide are linked to the content of Annex and Supplementary document, which serves as a more practical reference. In particular, Process 2, which requires a higher level of expertise, is emphasized and explained in detail.

## Guide to Red Teaming Methodology on AI Safety





Red teaming was conducted on an LLM system using RAG, and based on the results, the implementation points and examples of deliverables for each process are documented in Annex and Supplementary document.

## [Methods of red teaming]

<b>Conditions of red teaming</b>	Period	Two months from January to February 2025
	Attack vector	Conduct red teaming remotely from company devices through the internal network to the AI system in the development environment, using the input prompt as the attack interface.
	Red teaming scope	Entire AI system
	Definition of attack success	Identification of vulnerabilities that directly cause harm to individuals or organizations, as well as those that may not cause direct harm but have potential concerns for being exploited in other attacks.
<b>Environment of red teaming</b>	Target AI system	Internal business data utilization chatbot service for enterprises using RAG
	User scope	Company employees
	Main provided functions	Employees can upload a complete set of business-related documents to a restricted-access container, enabling the generation of responses based on those documents. This allows employees to use the chatbot as an internal document search service, facilitating flexible information retrieval through a conversational UI.
	AI system configuration	Region: East US 2 Pricing Tier: Standard S0 Base Model: GPT-3.5-Turbo (1106) Rate Limits: 100,000 TPM / 600 RPM *All GPT models in Azure OpenAI Service can be selected and utilized freely. ■ RAG [Search] Azure AI Search (Standard S1) [Data Source] Azure Blob Storage (Standard) [Data Extraction Engine] Azure AI Document Intelligence (Standard S0)

# AISI

Japan AI Safety Institute