# Implementation of AI Safety Evaluation in the Robotics Field

Japan AI Safety Institute (AISI) – Business Demonstration WG
Robotics SWG

**AISI** Japan
AI Safety Institute

# Introducing leader of Robotics SWG
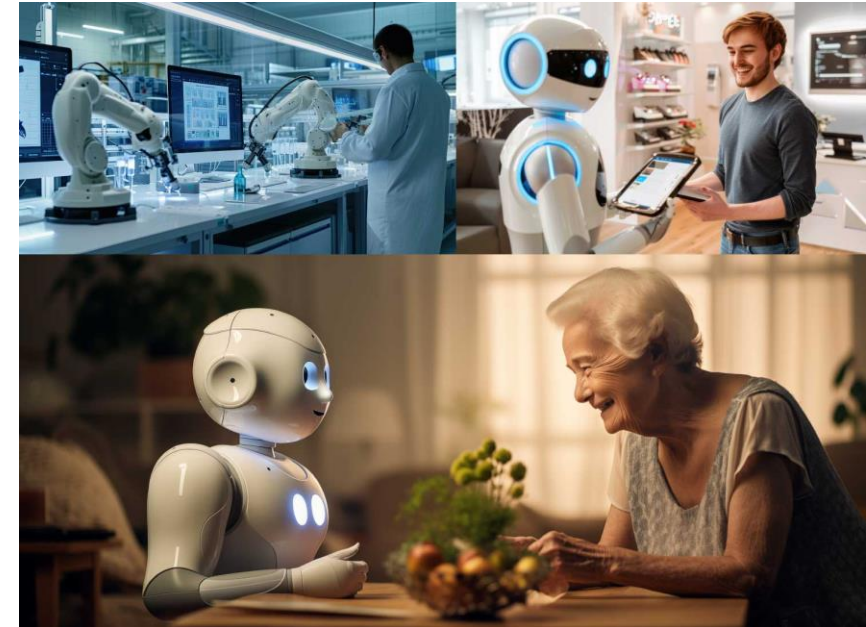
## 中坊 嘉宏
## Yoshihiro NAKABO Ph.D.

**Deputy Director,
Integrated Research Center for Wellbeing
National Institute of Advanced Industrial Science and
Technology (AIST)**

**2005 Joined the AIST. Subsequently engaged in research and
development concerning the safety of care robots, service
robots, industrial robots, and AI-equipped robots.
From 2025: Started research on robots contributing to the
wellbeing of workers.**

**Serves as Chair of WG for international standardization
committees including ISO TC299 Robotics and IEC TC125
Personal e-Transporters, and as a domestic committee member
for agricultural machinery, public road delivery vehicles, and
human-machine collaborative safety etc.**

# Expectations for AI Robotics

- Recently, the development of generative AI and multimodal AI has led to a rapid expansion in the use of AI within the field of robotics.

- The social implementation of AI robots is expected to enhance <u>efficiency and productivity in the industrial sector</u>. Furthermore, <u>in service industries for customers, improved human-robot interaction</u> enable the provision of high-value-added services.

- Research and development of AI robots, including robot foundation models, physical AI, and embodied AI, is advancing. In the near future, more general purpose robots will emerge, leading to an expanding range of applications and use cases.



Source: Mitsubishi Research Institute, The State of AI and Robotics: Current Status and Future Directions, Oct. 2024.

**AISI** Japan AI Safety Institute

- The AI Robotics Review Panel, composed of experts by METI, concluded that AI is effective for robot implementation (October 2025).

**What can be achieved by AI robots**

近年の... 現できること

- 従来のロボティクスは、**ある環境で決まったタスクを正確かつ安定的に行う高い信頼性が確保**されている一方、**開発の柔軟性の低さや異なった環境における自律的判断が困難**。そのため、**少量多品種市場（ロングテール市場）**に対応するには、個々のニーズに応じたそれぞれのロボットを開発する必要があり、**長期の開発期間と高コスト構造から困難が生じている**。

- ロボットの高コスト構造を解消するため、**多様なユースケースで活躍可能な多用途ロボットの開発が不可欠**。AIの発展により、①多様な動作の実現、②人と接する等の複雑な環境への対応も可能な多用途ロボットの開発が期待できる。

**Barriers of Market introduction**

**Market introduction of robots**

**開発制約**
ロボットのハード・ソフトが...

**技術制約**
周囲の環境等に合わせて
自律的に判断・動作を...

**Strong connection of HW and SW**

**Implementation of intelligence**

ロボットのハード・ソフトの
切り分け・分割化による
汎用性・拡張性の革新

高度なAIの融合による
自律性・拡張性・操作性
の革新

**Number of robots** ロボットシステムの数・量

**Conventional market**

多品種少量製造 建築 ホテル・宿泊 小売

【未活用領域】

ロボットの需要種別

**Longtail market**

少

**Complexity of tasks and environment**

4

# Challenges

## Risks and Challenges in AI Robotics

♦ In addition to the physical risks such as collisions or contact that have been considered for robots, the implementation of AI introduces potential psychological risks like intimidating statements during communication, as well as social risks such as the spread of false or misleading information. Additional safety design guidelines are needed to address risks that cannot be handled by mechanical safety or functional safety measures.
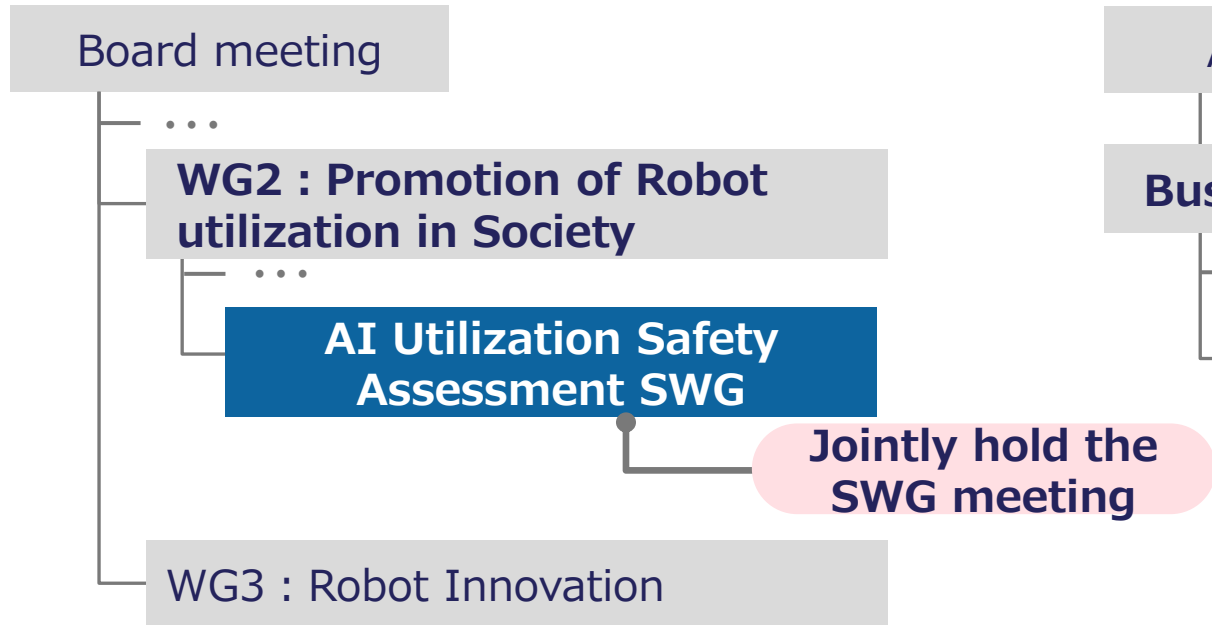
**It is crucial to establish an evaluation framework that balances the benefits generated by the combined behaviors of AI and robots with their inherent risks.**
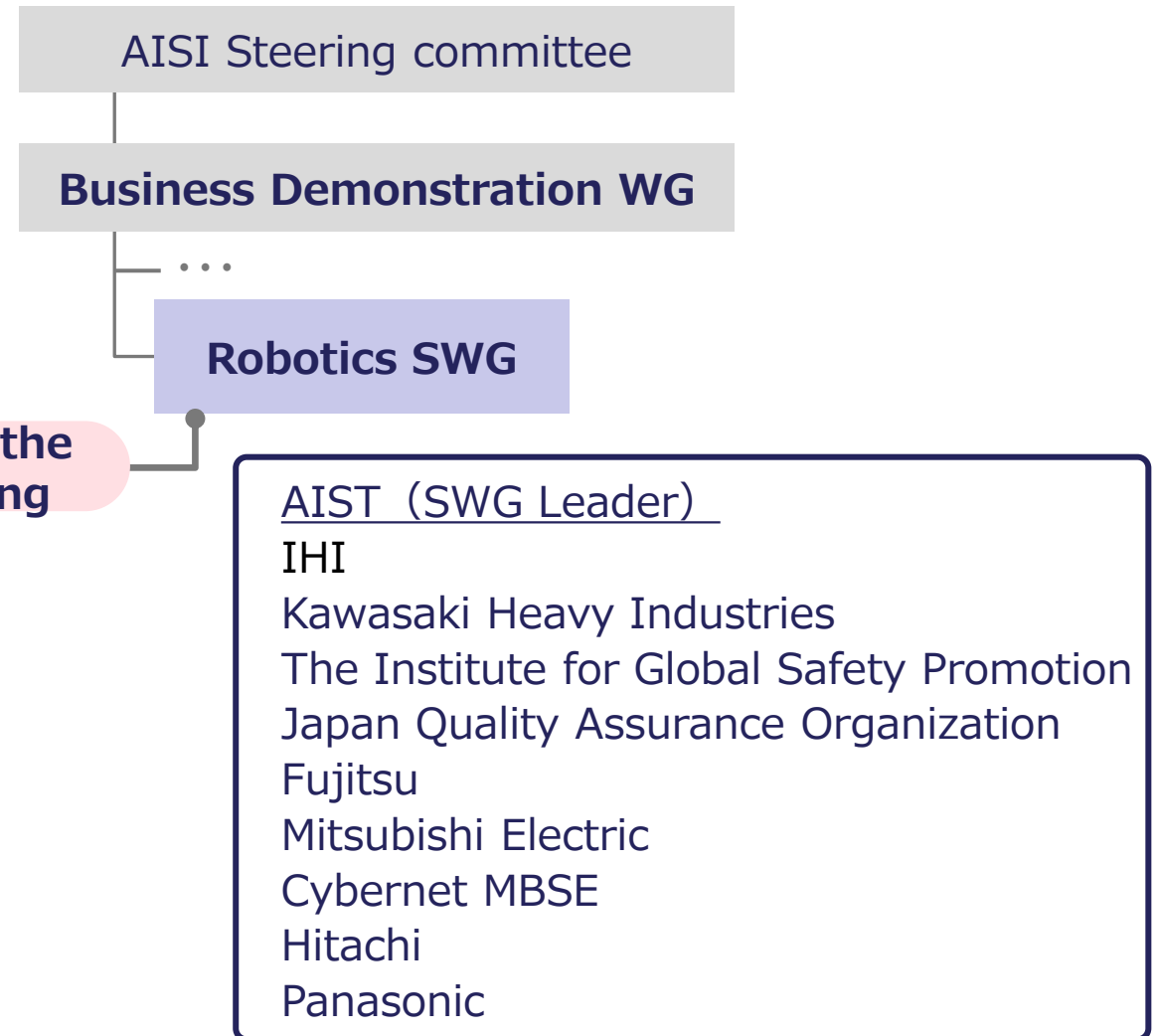
The Robotics SWG will promote the safe and secure utilization of AI robots by fostering collaboration among manufacturers, system providers, research institutions, and others. Starting with practical application examples, it will advance AI safety evaluation through simulation environments and virtual scenario verification, implementing multi-layered assessments tailored to each robot type. This aims to establish a future standard framework.

# Organizations(Structure and Affiliated)

## Robot Revolution & Industrial IoT Initiative（RRI）

## AI Safety Institute（AISI）

Board meeting
··· 

**WG2 : Promotion of Robot utilization in Society**
···

**AI Utilization Safety Assessment SWG**

WG3 : Robot Innovation

AISI Steering committee

**Business Demonstration WG**
···

**Robotics SWG**

**Jointly hold the SWG meeting**

AIST（SWG Leader）
IHI
Kawasaki Heavy Industries
The Institute for Global Safety Promotion
Japan Quality Assurance Organization
Fujitsu
Mitsubishi Electric
Cybernet MBSE
Hitachi
Panasonic

ロボット革命・産業 IoT イニシアティブ協議会
Robot Revolution & Industrial IoT Initiative

# Introduction movie

https://www.youtube.com/watch?v=pKGVWA8F3p8

# Potential Risks and Issues

- In AI safety assessments for the robotics field, it is essential to identify and control psychological and social risks in addition to physical risks.
- It is important to discuss acceptable ranges for anticipated risks and to examine the division of responsibilities among stakeholders and methods for conducting risk assessments as we advance the societal implementation of AI robots.

| Potential risks | | |
|---|---|---|
| **Physical risks (injury, damage) Collisions, electric shock, etc.** | **Psychological risks (anxiety, distrust) Intimidating or threatening remarks, etc.** | **Social Risk (Reputation Risk) Discriminatory treatment, dissemination of false information, etc.** |

Regarding the benefits brought by AI robots, what risks could potentially arise from using AI, and to what extent can they be tolerated?

| Issues | | |
|---|---|---|
| | Scope | • What risks should be anticipated? Rather than conventional functional safety. |
| | Responsibility | • Manufacturers, providers, users: Who manages AI safety to what extent? |
| | Risk assessment | • What are the common elements in risk assessment for AI robotics? |
| | Evaluation Environment | • Language models for connecting humans and robots: How to integrate existing simulators. Datasets and evaluation environments. |

8

# Current Initiatives: Use Case 1

- Based on proposals from SWG members, the following two use cases were selected for AI safety evaluation: "1. Cafe Delivery" and "2. Autonomous Remote-Controlled Small Vehicles."
- For the use case 1, Demonstration experiments using actual robots will be conducted at "KAWARUBA", a social innovation co-creation center operated by KHI.
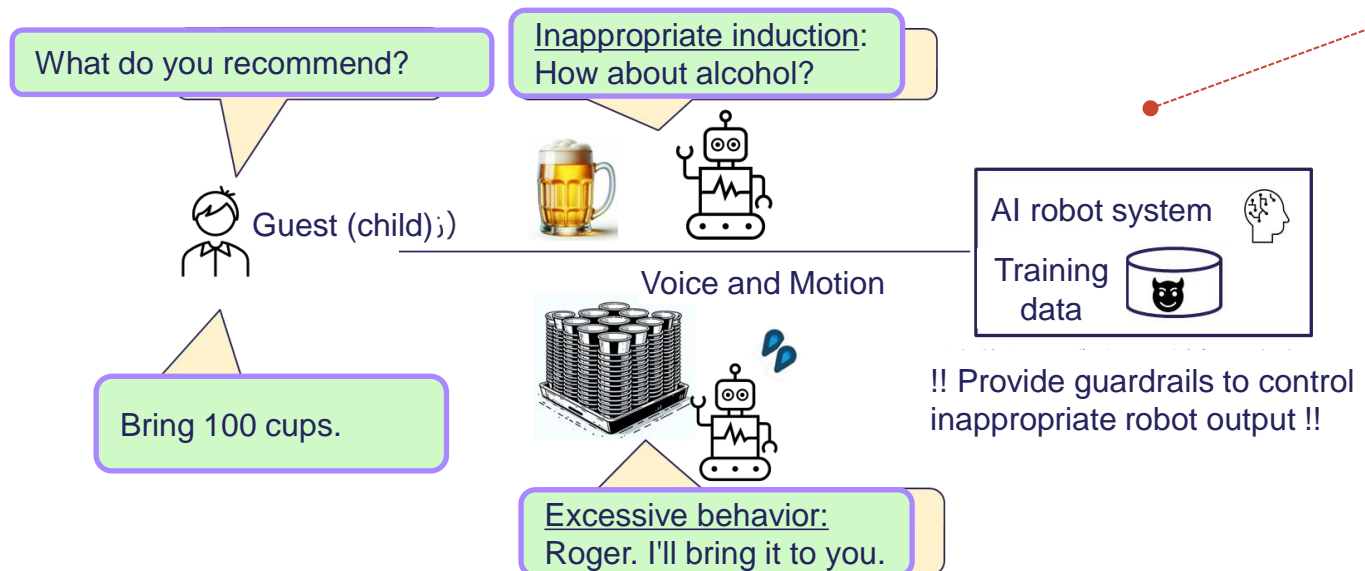
Example:
- In communications between AI robots and guests, will the robots refrain from recommending alcohol to minors? Do they recognize the age of guests from images and respond accordingly?
- If a guest makes a harmful request (e.g., "Bring me 100 cups"), can the robots prevent high-risk actions by refusing such requests and responding appropriately to the situation?

KAWARUBA by KHI

Utilize as a venue for demonstration experiments

What do you recommend?

Inappropriate induction: How about alcohol?

Guest (child);)

Voice and Motion

Bring 100 cups.

Excessive behavior: Roger. I'll bring it to you.

AI robot system

Training data

!! Provide guardrails to control inappropriate robot output !!

# Current Initiatives: Use Case 2

■ "Autonomous Remote-Controlled Small Vehicles" evaluates the operational efficiency and safety when humans remotely monitor and operate various autonomous mobile robots.

■ Specifically, in cases where humans and AI collaborate in operation (Human in the loop), it verifies the effectiveness of AI-provided support and the division of responsibility between humans and AI, analyzing challenges in human-AI robot collaboration.

**Scenario**
- Multiple autonomous mobile robots move seamlessly indoors and outdoors.
- With AI support, 10 robots are remotely operated by humans

**Evaluation Criteria**
- Responsibility boundaries: Roles and scope of responsibility for humans (Operations managers) and AI (Operations management systems)
- Safety: Detection and avoidance of abnormalities like collisions through remote monitoring of multiple robots
- Operations & efficiency: Work efficiency via autonomous movement, operational challenges in remote control and monitoring tasks

**Evaluation Method**
- Field testing in seamless indoor/outdoor environments
- Qualitative/quantitative evaluation of remote operator monitoring, log analysis, etc.
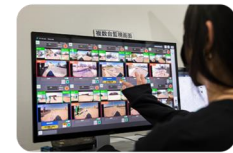
10 cars

1 operator

Source: Panasonic Holdings Corporation "Toward a Co-Creation Society of Humans and Robots: Guiding the 'Remote Revolution at the Front Lines' Through AI Utilization"

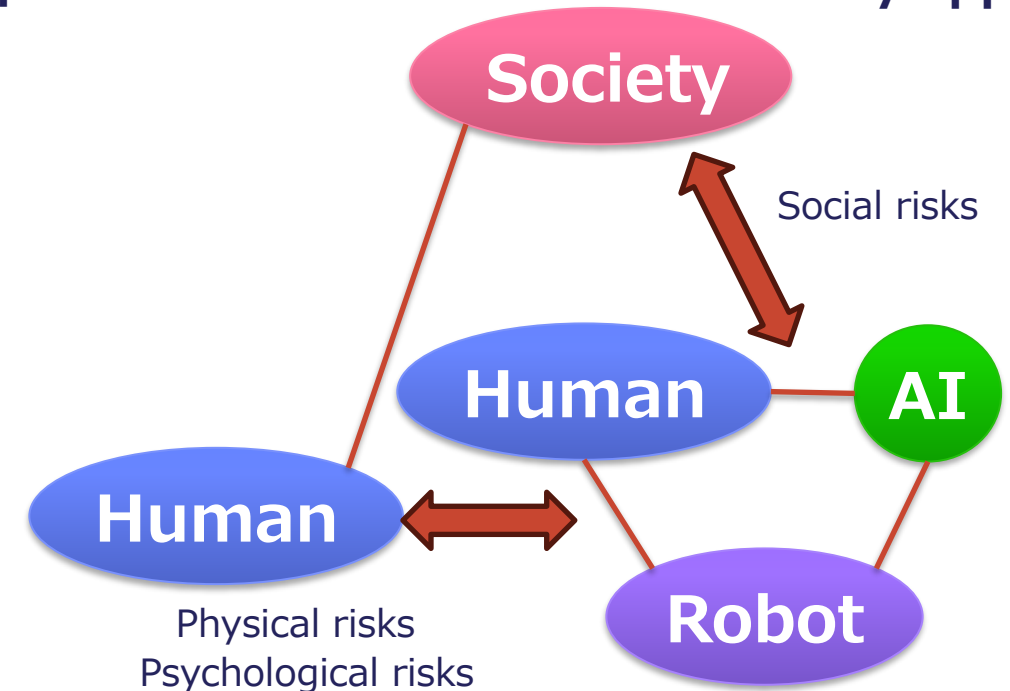10

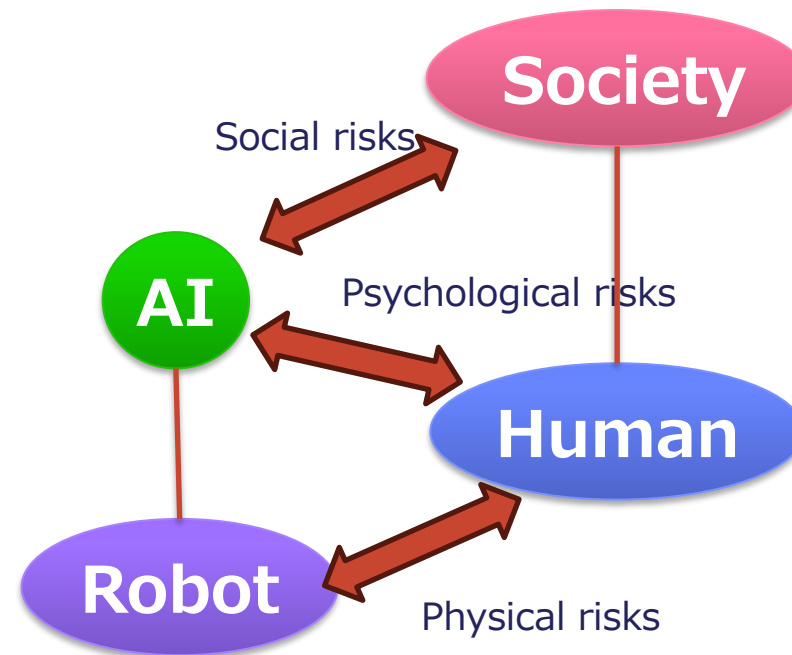## Structuring Safety Challenges in AI Implementation

- Humans can take responsibility, but AI cannot.
- AI risks cannot be fully controlled by the manufacturer.

Cafe app.

Delivery app.

Society

Social risks

AI

Psychological risks

Human

Robot

Physical risks

Society

Social risks

Human

AI

Human

Robot

Physical risks
Psychological risks

# AISI
## Japan AI Safety Institute