# HITACHI
## Inspire the Next

**EU-JAPAN DIGITAL WEEK 2025**

**"Critical Applications of AI in Industry, Healthcare and Other Sectors" Workshop**

# OT x AI safety approaches in Hitachi

April 7th, 2025

**Satoshi Otsuka**

Autonomous Control Research Department

Mobility and Automation Innovation Center

Research & Development Group

Hitachi, Ltd.

# Contents

1. Background: Hitachi/AI Safety Architecture
2. AI Safety Shell (AISS)
3. Safety Engineering for AI (evolving system)
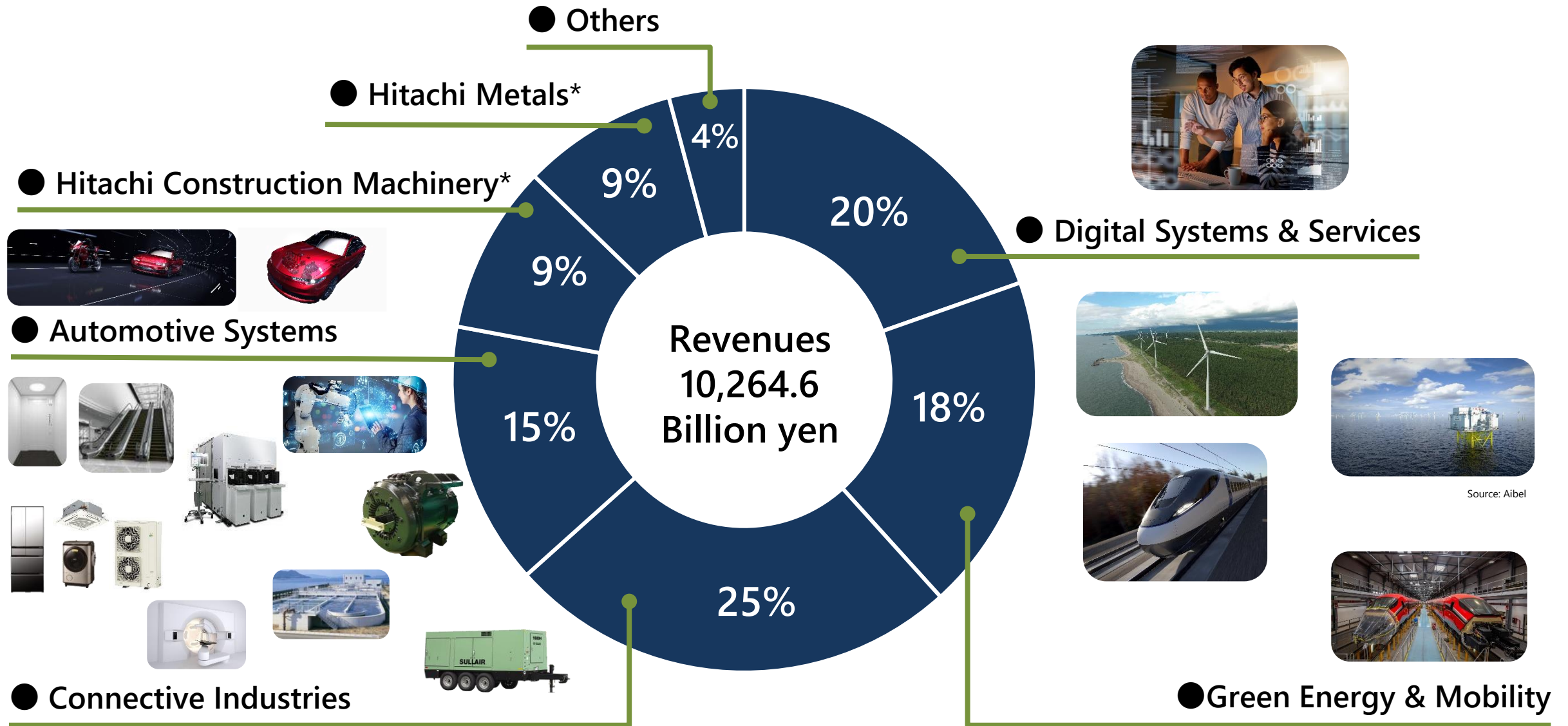4. Future Plan
5. Conclusion

# Contents

# 1-1. Hitachi corporate data

| Corporate Name | Hitachi, Ltd. |
|---|---|
| Founded | 1910 |
| Headquarters | 6-6, Marunouchi 1-chome, Chiyoda-ku, Tokyo 100-8280, Japan |
| Revenues | 10,264.6 billion yen (FY2021[*1]) |
| Adjusted operating income | 738.2 billion yen (FY2021[*1]) |
| EBIT (Earnings before interest and taxes) | 850.9 billion yen (FY2021[*1]) |
| Net income attributable to Hitachi, Ltd. stockholders | 583.4 billion yen (FY2021[*1]) |
| Number of consolidated employees | 368,247 (As of end of FY2021[*1]) |

*1: Based on the financial results for FY2020 ended in March 2021

# 1-2. Business segment constitution (FY2021)

● Others

● Hitachi Metals*

● Hitachi Construction Machinery*

● Automotive Systems

● Connective Industries

● Digital Systems & Services

● Green Energy & Mobility

**Revenues 10,264.6 Billion yen**
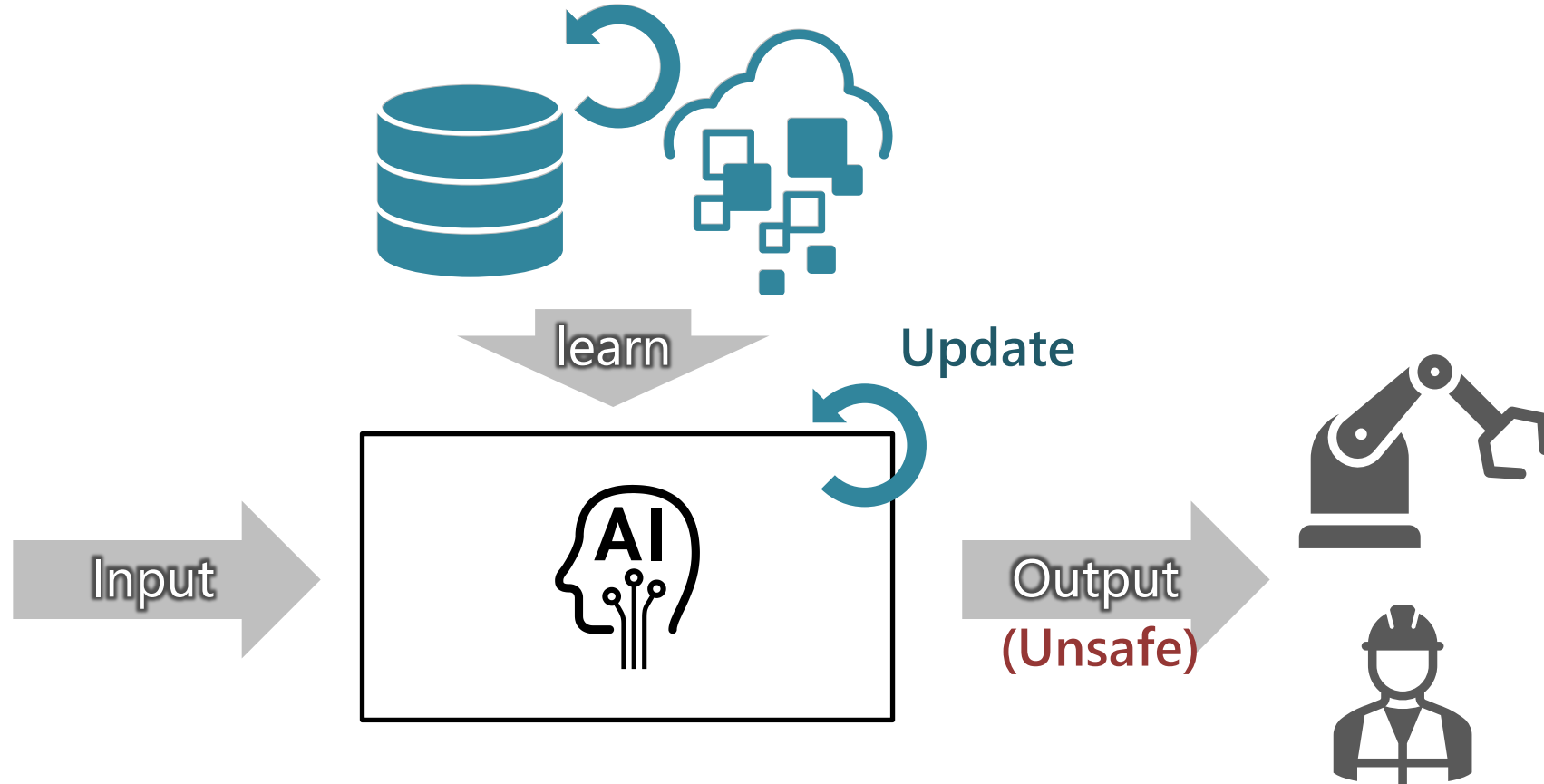
- 4%
- 9%
- 9%
- 20%
- 15%
- 18%
- 25%

Source: Aibel

The figures are based on the new segment classifications effective from FY2022.

*   Hitachi Construction Machinery was deconsolidated on August 23, 2022. Hitachi Metals are scheduled to be deconsolidated in FY2022.

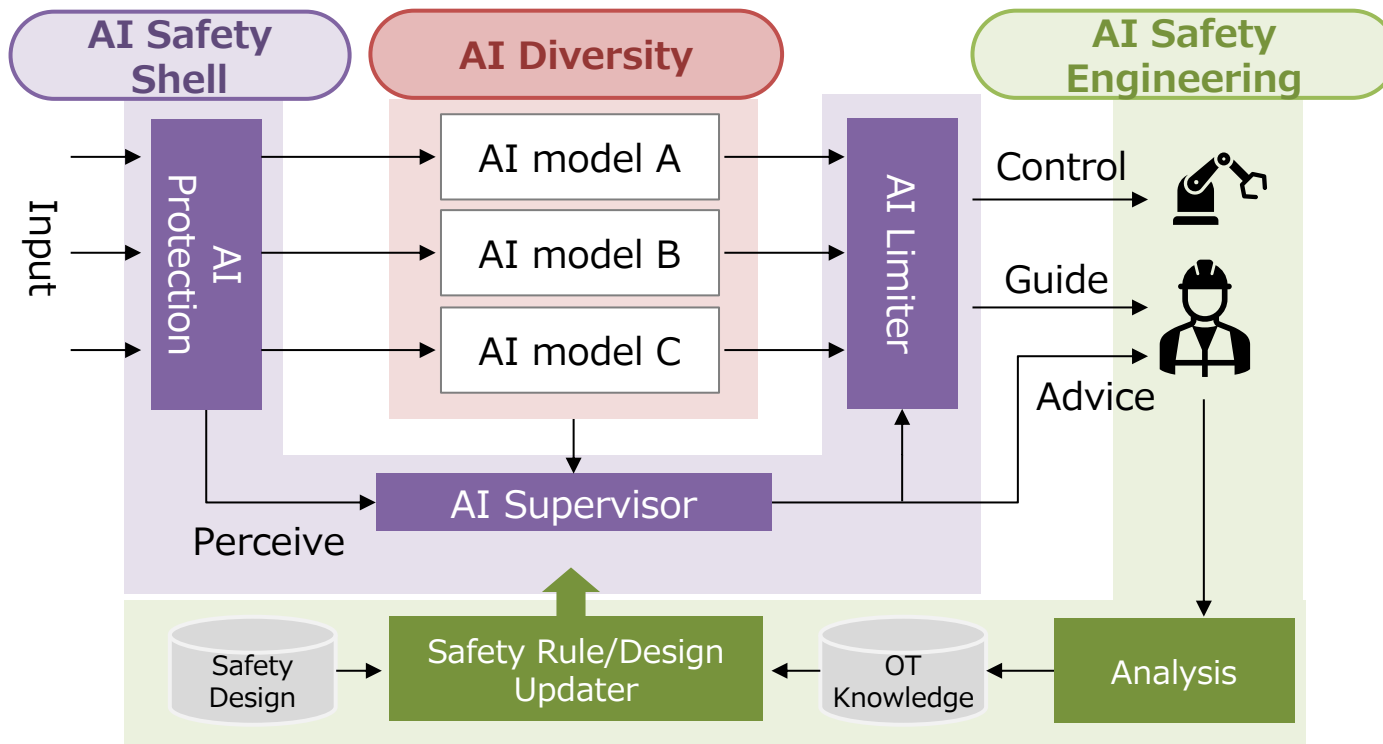**AI is a powerful tool. However, there are existing risks, and it cannot be used as it is.**



learn

Update

Input

AI

Output
(Unsafe)

# 1-4. Existing countermeasure methods (architecture)

HITACHI
Inspire the Next

**Architecture requirements: safe protection, no loss of efficiency, support for evolution**

| | Architecture Type (Classification) | Subtype | method | Pros | Cons |
|---|---|---|---|---|---|
| A | Making AI system more reliable | AI component protection | ・Guardrail [1] | ・Protection of unintended output is possible | ・Specific output can be suppressed, but unintended outputs are possible |
| B | | | ・Uncertainty separation [2] | ・Accuracy can be improved (Flip uncertain output) | ・Zero incorrect output is not possible (Uncertainty of AI itself) |
| C | | AI component redundancy | ・AI redundancy (Ensemble [3]) | (parallel processing) | |
| D | | | ・AI diversity (multi-agents) [4] | ・Accuracy can be improved (Select good results in parallel processing) | ・Zero incorrect output is not possible (Mistake by consultation) |
| E | Rule-based protection from the outside | Fixed limiter | ・Upper and lower limits/limit values [5] | ・Almost reliable protection | ・Performance degradation due to exc |
| F | | Evolutionary support | ・MAPE-K [6] ・Dynamic risk assessment, etc. [7] | ・The limiter also evolve according to the situation. | ・Di evo |

**Good to use it to improve accuracy (high reliability).**

**Essentially, impossible to guarantee the output of AI itself => Insufficient for "OT x AI" safety measures**

**Safety assurance requires certain protection (rule-based)**

**The design needs to be able to dynamically change the rules and follow them "correctly" and "efficiently".**

**HITACHI**
**Inspire the Next**

## Improving reliability/safety through diversity, AI Safety Shell, and Safety Engineering



**①AI Diversity**
(Excluded from today's explanation)
- Improvement of reliability/performance

**②AISS**
- Ensuring safety by making decisions based on rules

**③Safety engineering**
- Responding to evolving systems

7

# Contents

**The key is supervision and input/output protection.**
**The way to protect it differs depending on the application/models/purpose.**

■ **AISS:**

1. **Protection by safety rules**
   (When it is easy to define)

2. **Risk-based protection**
   (for mission-critical systems)

**If the requirements for cooperation can be defined and AI based on safety rules, this approach is possible
The output of AI can be appropriately constrained by contract. Even if it is broken, AISS can protect the system**

○Structure

○Example

AI System

Spec Information
(Contract)

Monitor

AI Safety Shell

AI Spec Definition/
Analytics (OT×AI)

Safety monitoring/
Safe State Transition Control

Formulation of safety rules
(Domain Knowledge utilization)

AI Control system — Asp · Monitoring · Control — Bsp — Other systems (or person) — Bsp · Monitoring · Control — Asp
Rule generation

○AI performance information

$$A_{sp} = \left(a_{a\_min\_brake}, a_{a\_max\_accel}\right)$$
$$B_{sp} = \left(a_{b\_max\_brake}\right)$$

○ Safety rules

$$d_{min} = \left[v_a\rho + \frac{1}{2}a_{a\_max\_accel}\rho^2 + \frac{\left(v_a + \rho a_{a\_max\_accel}\right)^2}{2a_{a\_min\_brake}} - \frac{v_b^2}{2a_{b\_max\_brake}}\right]$$

# 2-3. Experimental result

**Even if the AI breaks the rules, it can be controlled safely**
**When updating, just update the rules**

## AISS off



## AISS on



*While it's easy to notice and stop when a vehicle suddenly accelerates, it can be difficult to notice when a vehicle sometimes doesn't slow down, and it can be too late.

**Protect outputs according to risk (for mission-critical systems)**
**The key is context. Protect outputs by correctly understanding the current situation.**



Created with historical/business insights

From the exhaustive combination, select and create at human discretion

"Hazardous events that should not happen" used to output protection

"Hazardous situations" are used for analysis

Cause Analysis Methods (FTA/FMEA, etc.)

Context Diagram

○Calculating the risk value

Calculating the cost of damage

Risk value is calculated by multiplying the **accuracy of the context** by the amount of damage listed on the left.

Estimate the probability of each event occurring

*Presentation Only

# Contents

**AI evolves through data-driven methods. Corresponding safety engineering methods are required**

■**Safety engineering**

- Dynamic risk assessment in response to evolution
- Methods for updating safety rules
- etc.

# 3-2. Safety engineering method overview

## Evolve the system to adapt human-evolution (co-evolutional control) with guaranteed safety

### Symbiotic safety (basis) [8]

Exclusive area

Exclusive area

Area exclusive control

Fraunhofer IKS

### Dynamic risk and capability assessment [9]

Fraunhofer IESE



Risk area extents

| $width = w_c$ | $width = w_o$ |
|---|---|
| $height = h$ | $height = h$ |

Capability area

Risk area

Active signal from AMR to HW

HW notices AMR signal

Runtime Evidence

Runtime Evidence

Capability area extents

| $width = w_c$ | $width = w_0$ | $width = w_1$ |
|---|---|---|
| $height = h$ | $height = h$ | $height = h$ |

Assumption: No friction and max. speed

No steering failure diagnosed

Acceptable tire pressure

Acceptable tire profile

Acceptable illumination

Design time Evidence

### Paired safety rule structure [10]

Hazard

Safety Goal

General safety rules — SRg · · · Safety Req. (Not for cooperation) · · ·

Designed by SPL methods

Specific (top) safety rules — SRst · · · · · ·

Derived by safety analysis

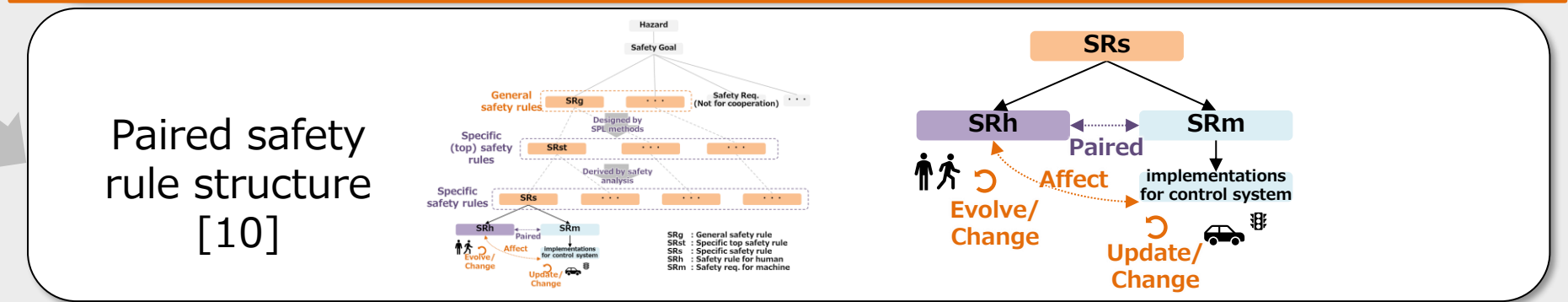Specific safety rules — SRs · · · · · ·

SRh Paired SRm
Evolve/ Affect implementations
Change for control system
Update/ Change

SRg : General safety rule
SRst : Specific top safety rule
SRs : Specific safety rule
SRh : Safety rule for human
SRm : Safety req. for machine

SRs

SRh ←Paired→ SRm

Evolve/ Affect implementations for control system
Change

Update/ Change

## Shorten the safety design time to apply the system to other use-cases or domains

### Environment metamodel [11]

Fraunhofer IESE

Safety engineering artifact → Reusable / New design

Argumentation + Safety engineering artifacts

DDI

16

# 3-3. Dynamic risk and capability assessment

**Dynamic Risk and Capability Areas Can Enhance Efficiency Compared to Using Worst-Case Assumptions**



- Replace static worst-case assumptions with dynamic safety reasoning capabilities
- Detecting low-risk situations and reducing risk/capability area extents
- Requires models that allow a dynamic assessment of:
  - risk of operational situation for human worker
  - safety-related capabilities of the AMR

AMR: Autonomous Mobile Robot

Engineering Runtime Safety Monitors Requires a Systematic Method for Dynamic Model Creation

SG: "Perform safe stop, when predicted dynamic risk and capability areas overlap"

Operational Context

Behavior Safety Concept

**Behavior Safety Analysis**

Actor Behavior Causality Models

Feature Selection & Formalization

Dyn. Risk Area Pred. Model

Dyn. Capability Area Pred. Model

Runtime Safety Monitor Component

Dynamic Risk Area for *Uncontrollable* Actors

Dynamic Capability Area for *Controllable* Actors

## Movement Space is Modelled based on SUDA and its Capability Decomposition

| Actor: AMR |
| --- |
| Movement Space Abstraction: Ellipse |
| Behavior: Move good safely through corridor according to plan |

| Actor: Human Worker |
| --- |
| Movement Space Abstraction: Ellipse |
| Behavior: Carry box from place A to B through corridor |

**Abstract Behavior Causality Model**

| Cognitive Step: Sense | Cognitive Step: Understand | Cognitive Step: Decide | Cognitive Step: Act |
| --- | --- | --- | --- |

| Capability: Pay attention to the scene | Capability: Localize other actors | Capability: Predict the actor behavior | Capability: Plan safe behavior | Capability: Execute planned behavior |
| --- | --- | --- | --- | --- |

**Behavior Causality Model AMR**

**Behavior Causality Model Human Worker**

Capability area

Risk area

*height*

*width*

*height*

*width*

SUDA: Sense-Understand-Decide-Act

## Potential Capability Deviations impact the likelihood of Behaviors

Capability area

*height*

*width*

*height*

*width*

Risk area

| Actor: AMR |
| Movement Space Abstraction: Ellipse |
| Behavior: Move good safely through corridor according to plan |

| Actor: Human Worker |
| Movement Space Abstraction: Ellipse |
| Behavior: Carry box from place A to B through corridor |

**Abstract Behavior Causality Model**

Cognitive Step: Sense → Cognitive Step: Understand → Cognitive Step: Decide → Cognitive Step: Act

Capability: Pay attention to the scene | Capability: Localize other actors | Capability: Predict the actor behavior | Capability: Plan safe behavior | Capability: Execute planned behavior

Deviation: No attention to scene | Deviation: Existing actors not localized | Deviation: Prediction critically wrong | Deviation: planned behavior is unsafe | Deviation: Safe plan executed unsafe

**Behavior Causality Model AMR** | **Behavior Causality Model Human Worker**

**Situation Features concretize abstract Deviations for particular Actors**

Lateral
Longitudinal

Capability area

$v_{HW}$

$v_{AMR}$

height

width

Occluded vision by packages

Risk area

**Actor: AMR**

Movement Space Abstraction: Ellipse

Behavior: Move good safely through corridor according to plan

**Actor: Human Worker**

Movement Space Abstraction: Ellipse

Behavior: Carry box from place A to B through corridor

**Abstract Behavior Causality Model**

Cognitive Step: Sense → Cognitive Step: Understand → Cognitive Step: Decide → Cognitive Step: Act

Capability: Pay attention to the scene

Capability: Localize other actors

Capability: Predict the actor behavior

Capability: Plan safe behavior

Capability: Execute planned behavior

Deviation: No attention to scene

Deviation: Existing actors not localized

Deviation: Prediction critically wrong

Deviation: planned behavior is unsafe

Deviation: Safe plan executed unsafe

DI: Cognitive Capacity

DI: Occlusion

DI: Intention Missinterpretation

DI: Perceived risk of situation

DI: Assumpt. of execution environment

SF: Lighting Condition

SF: Unobservable corner

SF: Prediction Model Assumption Validity

SF: FuSa fault leading to lower risk estimate

SF: Friction

**Behavior Causality Model AMR**

SF: Working Task

SF: Carried Box

SF: AMR Intent Indication

SF: Worker Experience

**Behavior Causality Model Human Worker**

DI: Deviation Influence,  SF: Situation features

21

**Integration of Dynamic Risk and Capability Models by ConSert**



## Passing through

Carrying a box

Not carrying a box & eye contact

Once the system detects worker's eye contact, Risk area is set to small

If the worker is carrying a box, risk area is set to large (he cannot see the AMR)

## Overtaking

If HW is not aware of AMR, AMR doesn't overtake

If HW is aware of AMR, AMR overtakes the HW

Once the system detects awareness, Risk area is set to small and thereby the AMR can overtake

22

# Contents

In this demonstration, we use a small-sized mobility instead of an autonomous forklift. Please watch the demonstration.

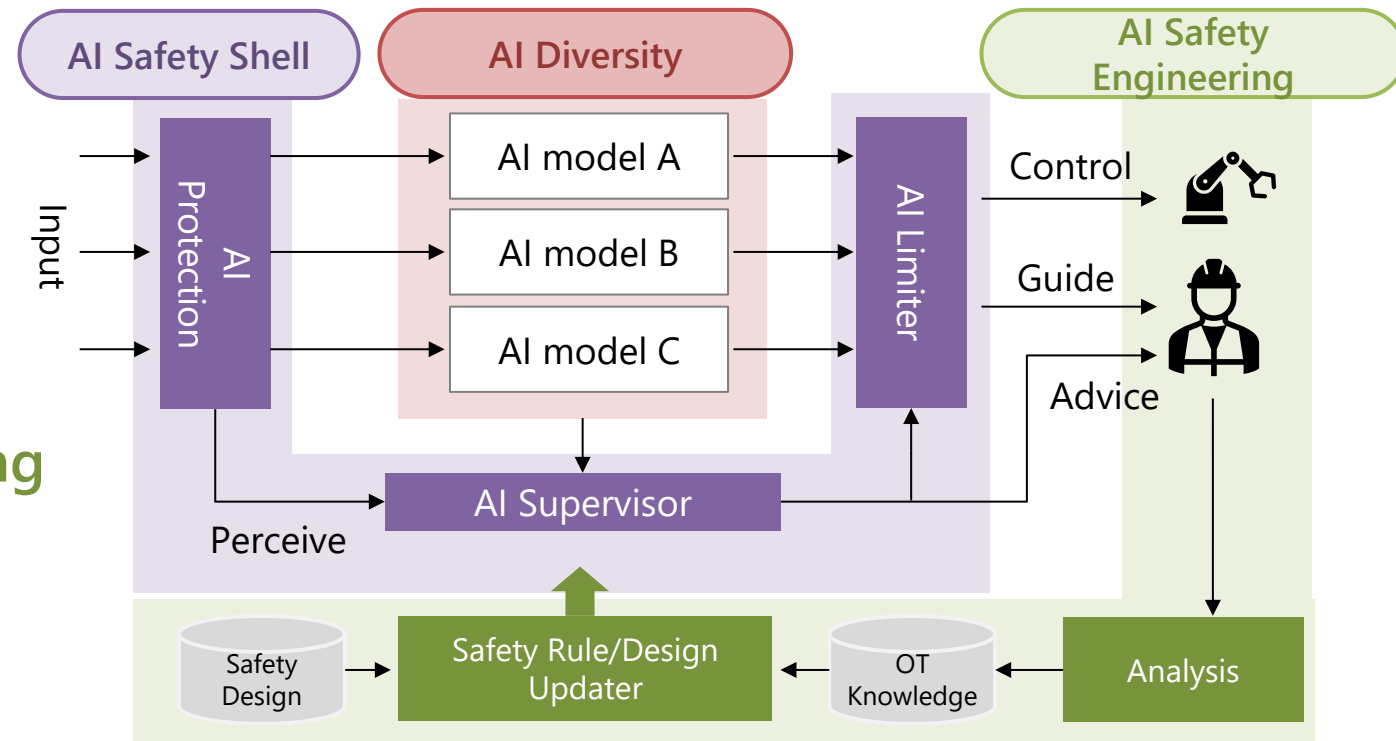**These contents are still being updated through the collaboration with Fraunhofer**

■**Specific implementation for AISS**
• Definition of AI risk
• Protection/Supervisor/Limiter

■**Updated version of Safety Engineering (by AI/ for AI)**

# Contents

- **Introduction of Hitachi's AI Safety activities**
  - Hitachi try to apply AI to Mission-critical systems (in the OT field)

- **Implementation of the AISS concept/architecture**
  - Two approaches: rule-based, risk-based assessment

- **Safety engineering process** that can respond to dynamic changes

- **Further activities**
  - Evaluation of implementation in control systems
  - Collaboration with Fraunhofer

[1] Dong, Yi, et al. "Building Guardrails for Large Language Models." arXiv preprint arXiv:2402.01822 (2024).

[2] Kläs, Michael, and Lena Sembach. "Uncertainty wrappers for data-driven models: increase the transparency of ai/ml-based models through enrichment with dependable situation-aware uncertainty estimates." Computer Safety, Reliability, and Security: SAFECOMP 2019 Workshops, ASSURE, DECSoS, SASSUR, STRIVE, and WAISE, Turku, Finland, September 10, 2019, Proceedings 38. Springer International Publishing, 2019.

[3] Petrakova, Aleksandra, Michael Affenzeller, and Galina Merkurjeva. "Heterogeneous versus homogeneous machine learning ensembles." Information Technology and Management Science 18.1 (2015): 135-140.

[4] https://note.com/fladdict/n/n106b9ce8f7d4

[5] Workgroup, E. G. A. S. "Standardized E-gas monitoring concept for gasoline and diesel engine control units." Version 5 (2013): 38.

[6] Weiss, Gereon, et al. "Towards integrating undependable self-adaptive systems in safety-critical environments." Proceedings of the 13th International Conference on Software Engineering for Adaptive and Self-Managing Systems. 2018.

[7] Reich, Jan, and Mario Trapp. "SINADRA: towards a framework for assurable situation-aware dynamic risk assessment of autonomous vehicles." 2020 16th European dependable computing conference (EDCC). IEEE, 2020.

[8] Ishigooka, T., et. al., "Symbiotic Safety: Safe and Efficient Human-Machine Collaboration by utilizing Rules," Design, Automation & Test in Europe Conference & Exhibition (DATE), 2022.

[9] Reich, J., et. al., "Engineering Dynamic Risk and Capability Models to Improve Cooperation Efficiency Between Human Workers and Autonomous Mobile Robots in Shared Spaces," Model-Based Safety and Assessment, pp.237-251, 2022.

[10] Otsuka, S., et. al., "Paired Safety Rule Structure for Human-Machine Cooperation with Feature Update and Evolution," Computer Safety, Reliability, and Security SAFECOMP 2023 Workshops, pp.247-259, 2023.

[11] Reich, J., et. al., "Concept and metamodel to support cross-domain safety analysis for ODD expansion of autonomous systems," Computer Safety, Reliability, and Security, pp. 165-178, 2023.

Hitachi Social Innovation is
POWERING GOOD