# Transparency and Integrity in HAIP Reporting Framework
## - Insights and Recommendations for AISI Network -

Future Developments of HAIP: Initial Reporting Outcomes and Alignment with the Japan AI Act:
Side event to the AISI International Network Meeting, Vancouver, Canada

18 July 2025

Arisa Ema (The University of Tokyo / Japan AISI)

Fumiko Kudo (The University of Osaka)
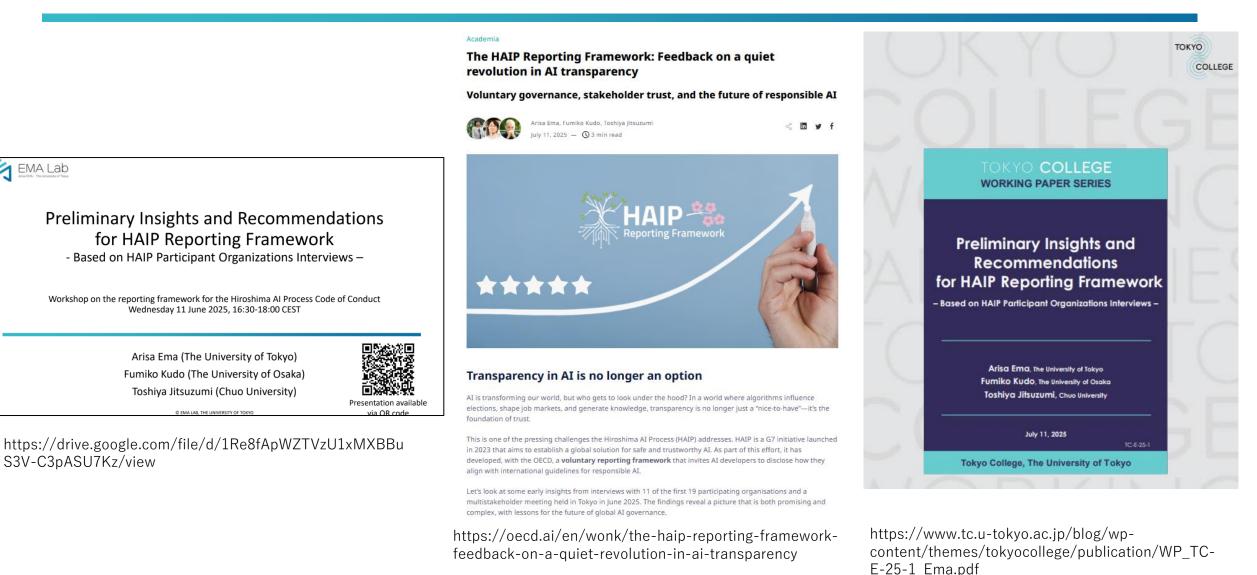
Toshiya Jitsuzumi (Chuo University)

Presentation available
via QR code

# Interviews with HAIP Participant Companies



https://drive.google.com/file/d/1Re8fApWZTVzU1xMXBBu
S3V-C3pASU7Kz/view



https://oecd.ai/en/wonk/the-haip-reporting-framework-
feedback-on-a-quiet-revolution-in-ai-transparency



https://www.tc.u-tokyo.ac.jp/blog/wp-
content/themes/tokyocollege/publication/WP_TC-
E-25-1_Ema.pdf

# Who Are the Target Audiences for HAIP Reporting?

| Audience Type | Description | Typical Motivation |
|---|---|---|
| **International Bodies** | G7 / OECD Partners | - Visibility in AI governance<br>- International alignment |
| **Policy Stakeholders** | Government bodies, regulators | - Gain trust<br>- Influence on regulatory frameworks |
| **Business & Technical Partners** | B2B clients, external developers, corporate partners | - Contractual clarity<br>- Risk accountability |
| **General Public** | Shareholders, citizens, job-seeking students | - Trust-building<br>- Brand strategy |
| **Internal teams** | Employees | - Create internal alignment and awareness on AI governance |

# What Effort Did HAIP Participation Require from Organizations?

- Reorganization of existing info vs. creation of new materials
  - Internal practices were sometimes not documented, structured for external audiences

- Internal approval hurdles (especially for Japanese companies)
  - Convincing internal teams of why transparency reporting matters
  - The submission deadline coinciding with fiscal year-end in March (in Japan)
  - Desire for broader understanding of HAIP's purpose and brand

# Ambiguities & Misunderstandings in the HAIP Questionnaires

- Ambiguities in:
  - Scope: Is the question referring to a specific AI system or company policies?
  - Role: Should we answer as a developer, a provider or both?
    - In B2Bcases, disclosure to clients can be particularly sensitive or difficult.
  - Audience: Is the report for government, clients, or the public?

- Needed for clearer templates or examples
  - However, there are tension between flexibility and clarity

# Report should Promote Transparency – Not Scoring

- Many companies assert that HAIP reporting framework should not be used for ranking/scoring without considering business model differences

  - Submitting a report demonstrates a commitment to transparency and responsible AI — this act itself should be encouraged

  - Integrity matters – need to prohibit unfair or deceptive acts or practices
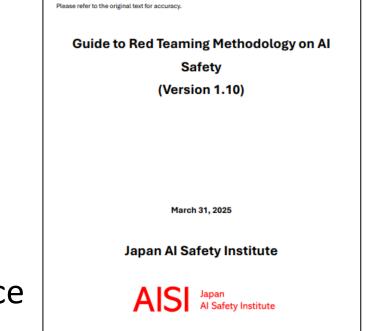
# Integrity matters – pre and post

- Entities responsible for ensuring the integrity of the HAIP report

| Level | Function | Actor |
|---|---|---|
| **Expert Guidance, Advice and Support (pre and post)** | Help companies write accurate report | AISI network / OECD-GPAI / UN |
| **Oversight and Monitoring** | Prohibit unfair or deceptive acts or practices | Government agencies / Authorities / Courts |
| **Social Accountability** | Detect and deter false claims and raise literacy | Civil society / Market / Journalism / Academia |

# Expected Roles for AISI in Supporting HAIP

- **Technical Expert** advisory
  - Identify and disseminate good practice
  - Provide pre- and post-report advisory sessions
  - Offer templates and tutorials
  - Serve as trusted consultation point

- Japan AISI already published Red Teaming Guidance
  - It could extend to report writing support!

- AISI network could help disseminate good practices and guide companies in choosing and explaining their approaches with integrity
  - This will be a support for SMEs as well

Please refer to the original text for accuracy.

**Guide to Red Teaming Methodology on AI Safety**

**(Version 1.10)**

March 31, 2025

**Japan AI Safety Institute**

**AISI** Japan AI Safety Institute

https://aisi.go.jp/output/output_framework/guide_to_red_teaming_methodology_on_ai_safety/

| HAIP Section | Key Technical Aspects required |
|---|---|
| **1 Risk Identification and Evaluation** | • Conducting **technical testing** (e.g. red-teaming, penetration test) to assess AI system readiness before deployment<br>• Identifying vulnerabilities, misuse through **adversarial testing** |
| **2 Risk Management and Information Security** | • Performing **testing** in secure, isolated or sandboxed environments<br>• Implementing robust **cybersecurity risk assessments**<br>• Protecting proprietary AI elements (e.g., model weights, algorithms) through **access controls and encryption** |
| **3 Transparency Reporting on Advanced AI Systems** | • Publicly disclosing **detailed results of technical evaluations**<br>• Providing information on **model capabilities, limitations, and appropriate use domains** derived from technical assessments |
| **4 Organizational Governance, Incident Management and Transparency** | • - |
| **5 Content Authentication & Provenance Mechanisms** | • Developing and implementing **technical mechanisms** (e.g., watermarking, metadata tagging, digital signatures) to identify AI-generated content<br>• Adhering to **international technical standards and best practices** for content provenance |
| **6 Research & Investment for AI Safety and Risk Mitigation** | • Investing in and conducting **research to develop new technical evaluation methods and tools** for AI safety, security, and trustworthiness.<br>• Advancing research in areas like **bias detection, disinformation, robustness, and explainability** through technical means |
| **7 Advancing Human and Global Interests** | • - |

# Next steps

- Shared Goals
  - Promote transparency in AI governance
  - Improve comparability across reports
  - Preserve flexibility and adaptability for diverse actors

- Next stems
  - Our detailed report and recommendations will be compiled this summer
  - We welcome feedback and continued dialogue from all stakeholders

# Special thanks

We sincerely thank the following organizations and individuals for their cooperation in the interview process:

**Organizations (by submission order):**

KDDI Corporation, SoftBank Corp., Preferred Networks, NEC Corporation, NTT, Microsoft, Salesforce, Anthropic, OpenAI, Google, Fujitsu, Rakuten Group

Additional organizations were invited, and we look forward to including their input in future versions.

Presentation available via QR code