

AI for Good

Global Summit

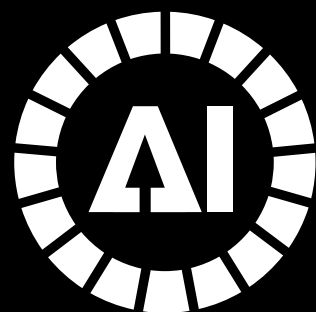
Introduction Slides

Mr. Kai WEI

Director, CAICT AI Institute

11 July 2025
Geneva, Switzerland





AI for Good
Global Summit

Characteristics of AI agents

CAICT

Agent

Agent

.....

Agent

LLM ~ OS

Perception

recognition,
comprehension,
interaction and
reasoning

Planning

planning, scheduling
and optimization

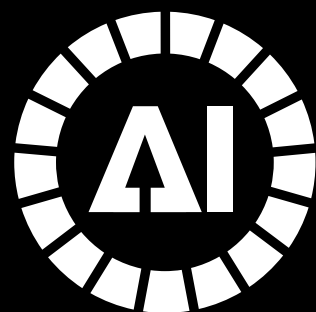
Action

execution in the digital
space or in the physical
environment

Memory

short-term and long-term
memory





Challenges

Reliability



- The root cause of agent reliability issues lies in **model hallucination**.
- To relief of the **hallucination is very hard**.

Safety & Security



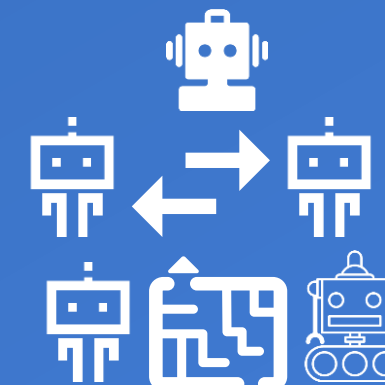
- Confronted with **prompt injection** and other emerging **attacks**.
- Traditional security risks such as **sensitive personal information** and privacy breaches are intensifying.

Interoperability



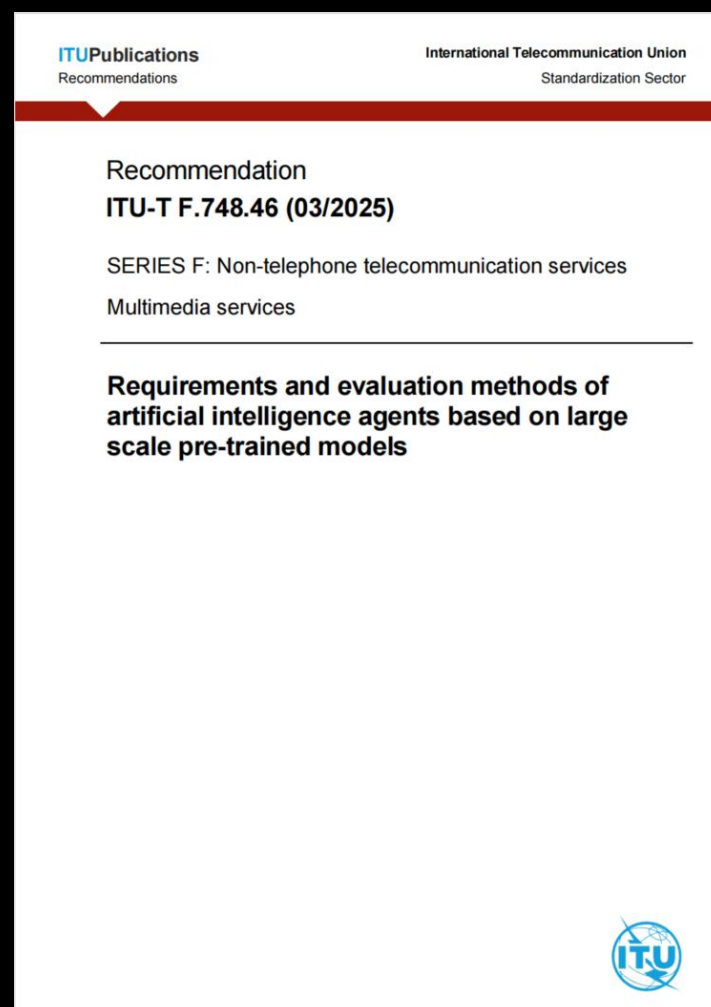
- Agent communication requires unified **interface standards**.
- The development of MCP and A2A ecosystems still requires **standards** for interoperability.

Operations & Maintenance



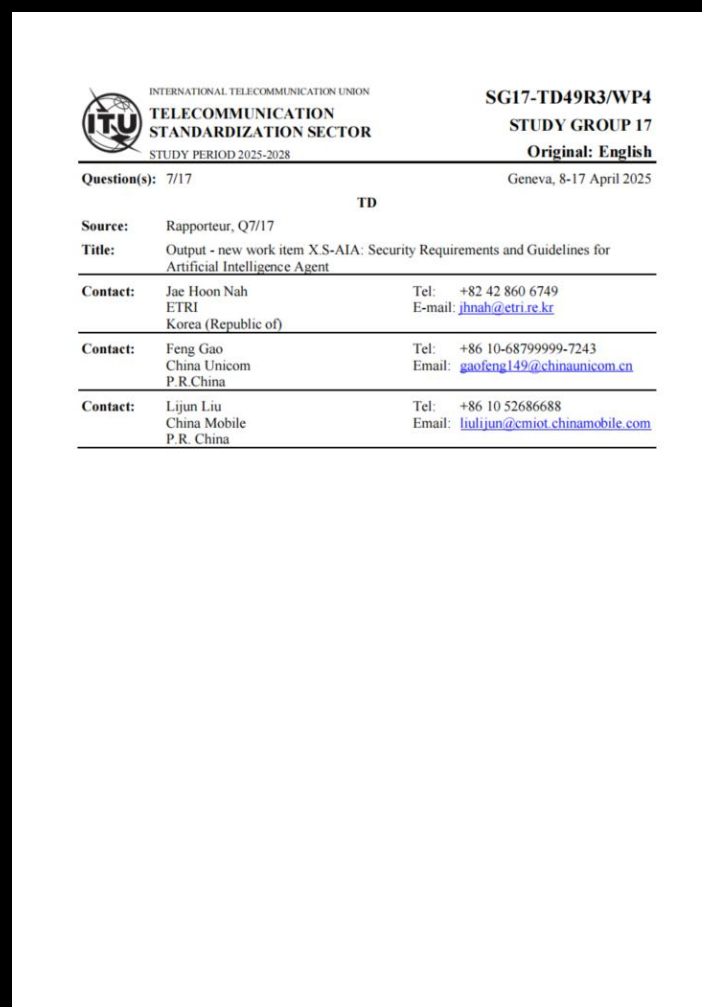
- The interactions between multiple agents increased the system **complexity**.
- It will significantly decrease the observability and increase the **difficulty of O&M**.

ITU-T F.748.46 @ SG21



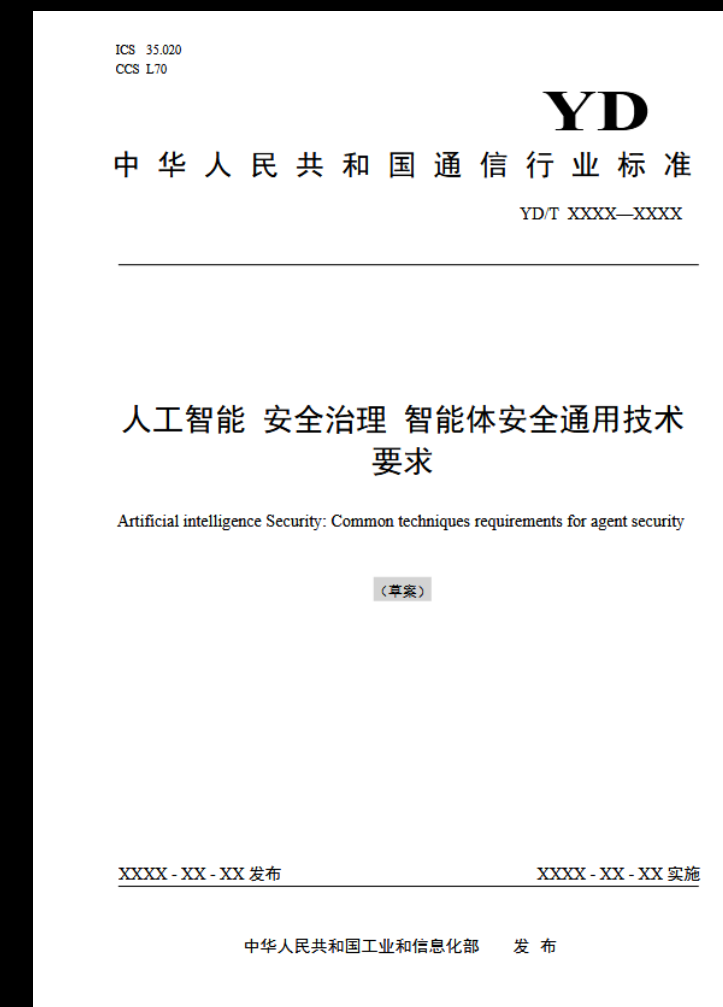
- Define the **capabilities and evaluation methods** for general AI agents.
- Propose an evaluation baseline for **agent-based product and applications**.

Security Req. of AI Agent @ ITU-T SG17

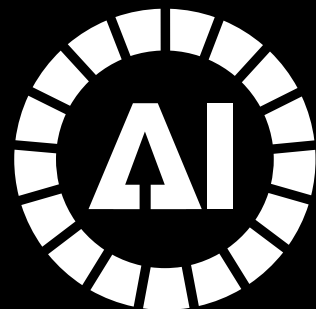


- Analyze security risks systematically across an agent's **perception, planning, memory and action** stages.
- Propose **lifecycle-wide risk protection** requirements.

General Req. of AI Agent security@ MIIT TC1



- Define the general requirements for **security capability of AI agent**.
- Establish a framework for **system interaction security** specifications.

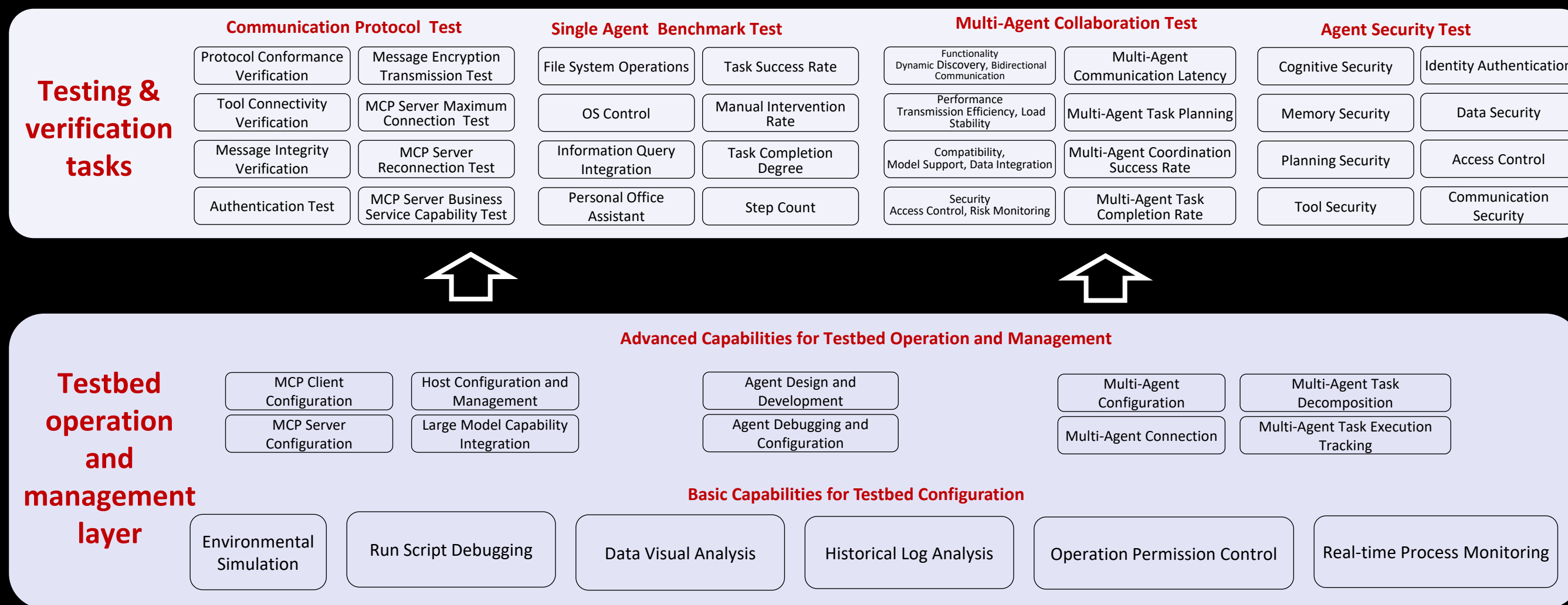


AI for Good
Global Summit

CAICT's Trusted AI Agent Testbed



CAICT is building a comprehensive AI Agent testbed for testing & verification of reliability, interoperability and security.



Source: Framework of Trusted AI Agent Testbed, CAICT, June 2025

