



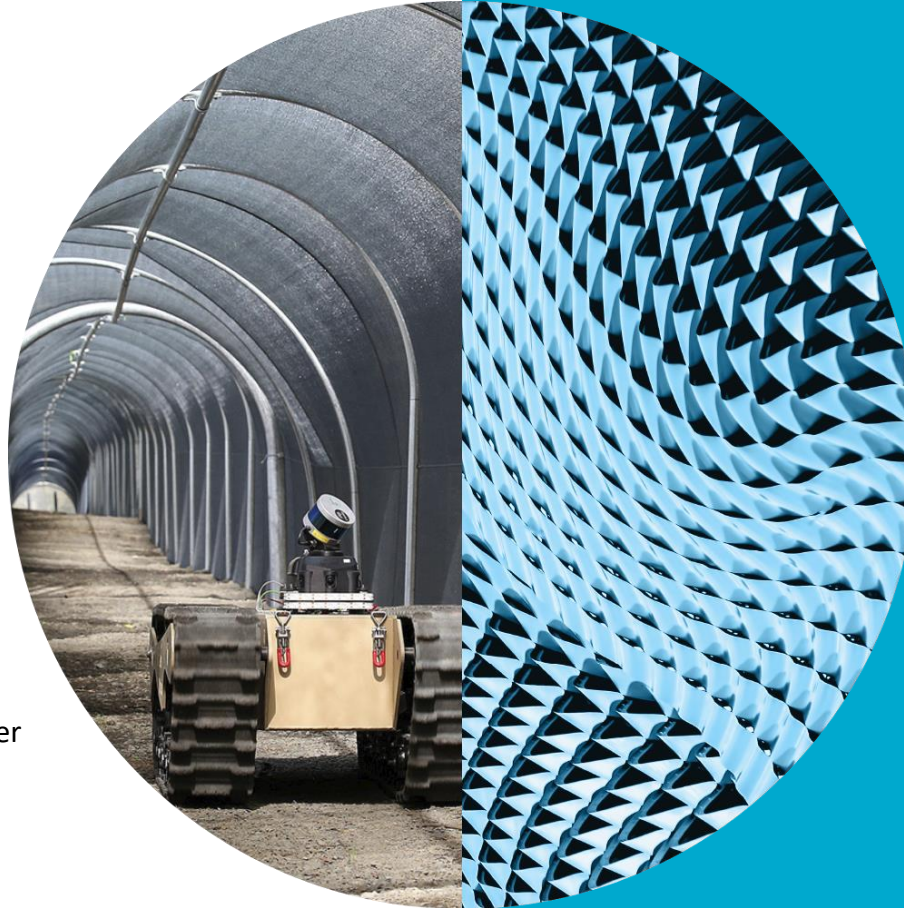
# Safe and Responsible Always Engineering

Qinghua Lu

Responsible AI Science Team Leader  
Data61, CSIRO

[qinghua.lu@data61.csiro.au](mailto:qinghua.lu@data61.csiro.au)

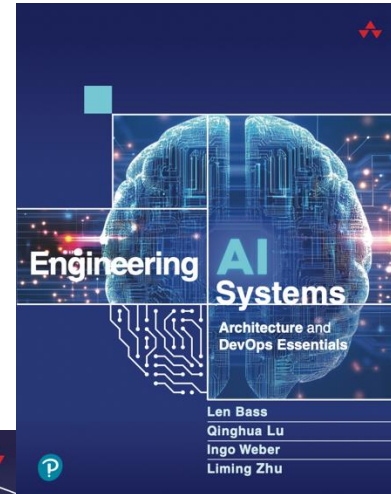
Australia's National Science Agency





# CSIRO and Responsible AI Team

- **CSIRO**
  - Australia's national science agency
  - Formed in 1916
  - 5500 people
  - 50 sites (Australia, France, Chile, US)
  - Data61: Data and Digital RU
- **Responsible AI team**
  - Formed in 2022
  - ~30 full time research scientists/engineers
  - Diverse and multidisciplinary team
  - 6 scientists in the top 30 for Responsible AI





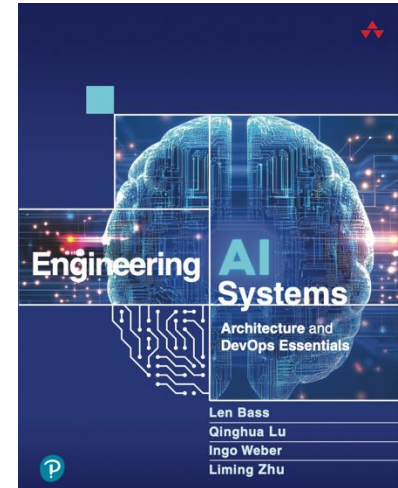
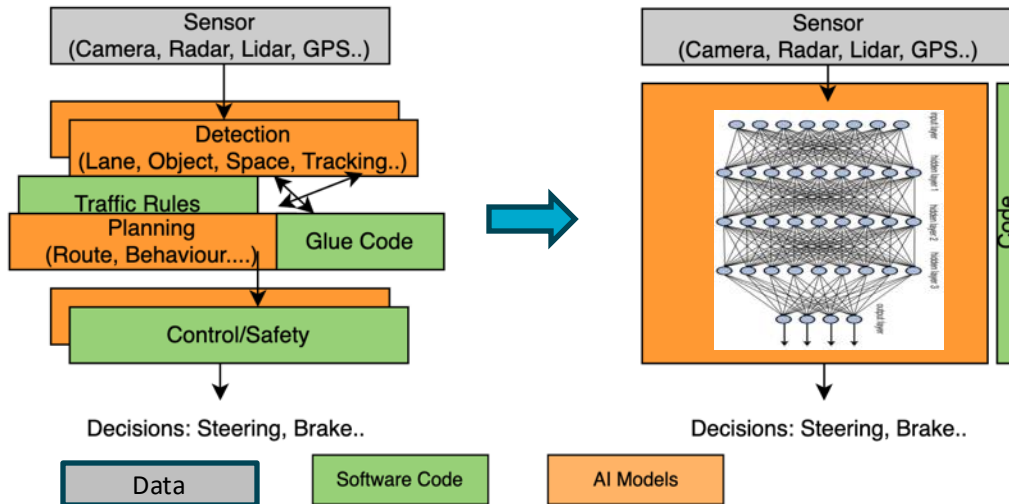
# What is an AI system?

- *An AI system is a machine-based system that, for explicit or implicit objectives, infers, from the input it receives, how to generate outputs such as predictions, content, recommendations, or decisions that can influence physical or virtual environments. Different AI systems vary in their level of autonomy and adaptiveness after deployment. (OECD, 2023)*
- *Artificial Intelligence (AI) is the research and development of mechanisms and applications of AI systems. (ISO 22989).*

# What is an Alware?

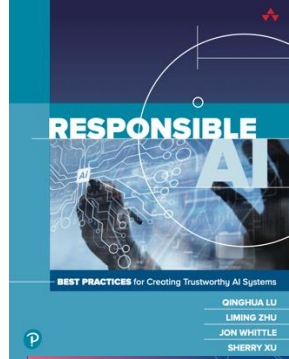
- *AI-as-Software, also known as Alware, refers to AI systems where functions are primarily encapsulated within a single general AI model as parameters/weights, rather than distinct narrow AI models explicitly chained together by traditional business code logic. (Bass, 2025)*

## End-to-End AI: Data In, Decision out, No Code



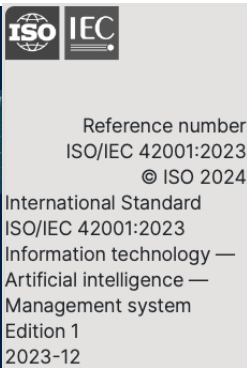
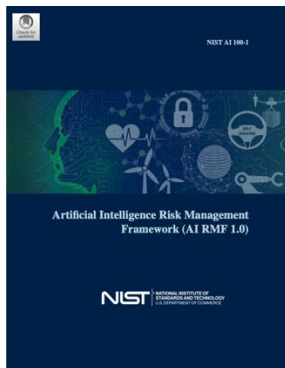
# What is Responsible/Safe AI?

- **Responsible AI** is the **practice of developing and using AI systems in a way that provides benefits to individuals, groups, and wider society, while minimizing the risk of negative consequences.** (Lu, 2023)
- **AI safety** is often used to describe **prevention of or protection against AI-related harms.** (Bengio, 2024)



Proposal for a Regulation of the European Parliament and of the Council Laying Down Harmonised Rules on Artificial Intelligence (Artificial Intelligence Act) and Amending Certain Union Legislative Acts  
2021/0106 (COD)

European Commission



### International network for AI safety

Evaluation of AI systems, foundational research, facilitation of information exchange

AISI	AISI	AISI	TBD	TBD
AISI	AISI	AISI	TBD	TBD

AI Office  
Regulatory body

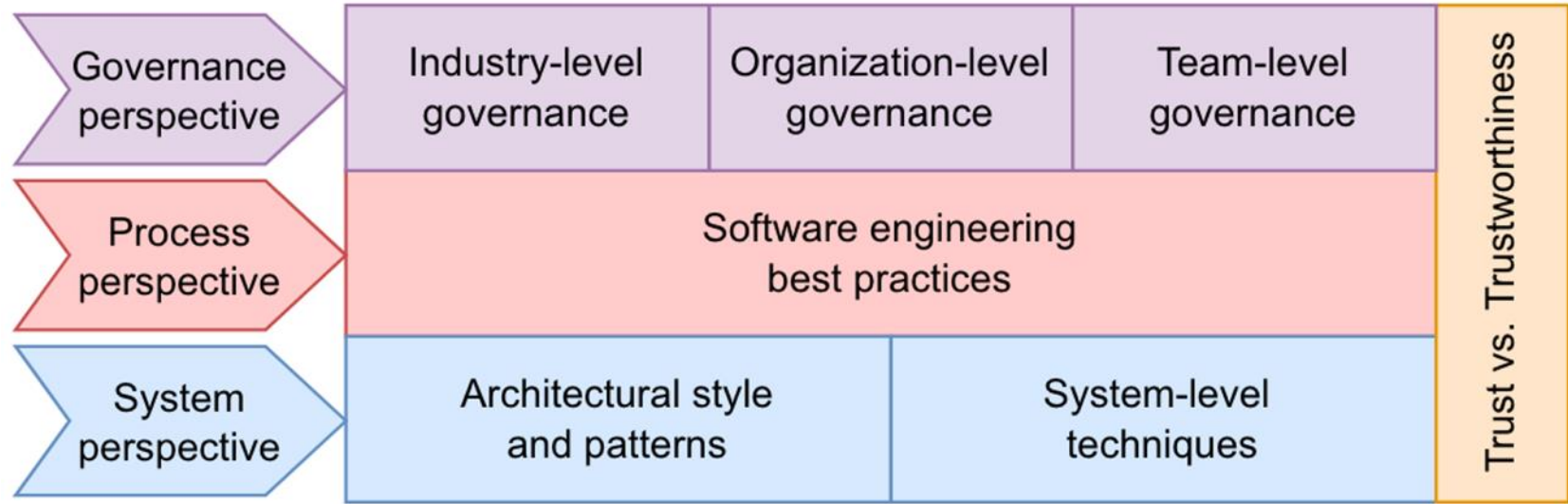




**Principles  
Standards**



**Responsible AI Engineering**

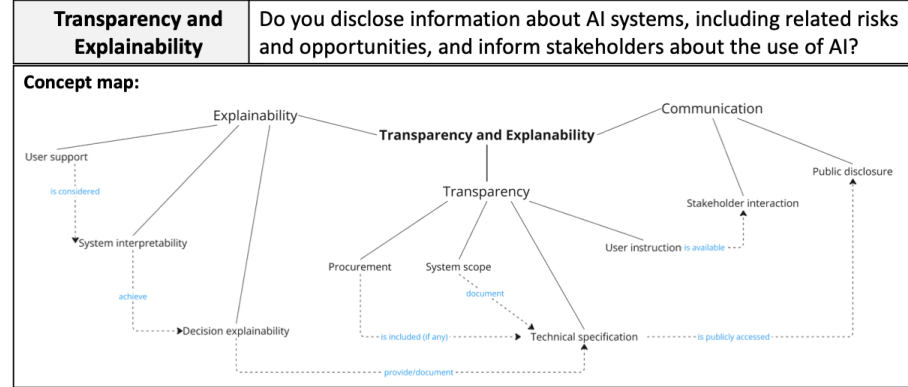
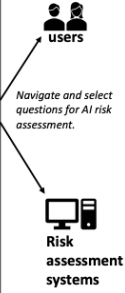
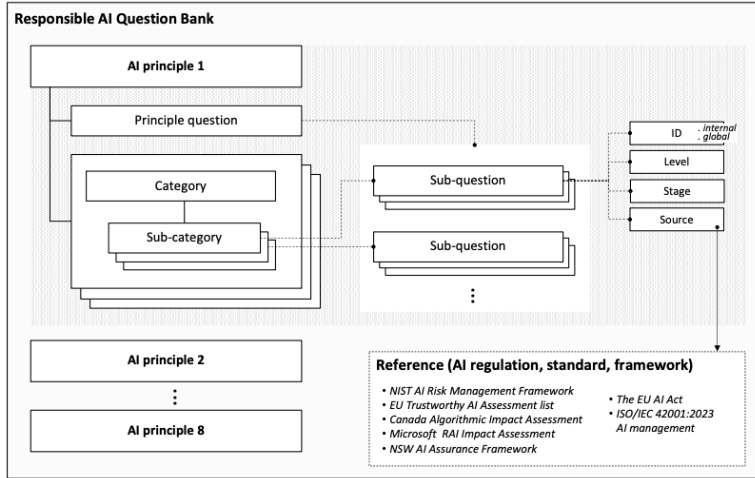


**Models**



Lu, Q., Zhu, L., Xu, X., Whittle, J., Xing, Z., 2022. Towards a Roadmap on Software Engineering for Responsible AI, in: 1st International Conference on AI Engineering (CAIN)

# Responsible AI Question Bank



**Sub-question (examples):**

Category		Question	Level	Stage	Source
Explainability	System interpretability	Do you design the AI system with interpretability in mind from the start?	2	D	EU
Transparency	Technical specification	Do you comprehensively understand and document the AI system including intended purposes, potentially beneficial uses, etc.?	2	P	NIST
Communication	Stakeholder interaction	Do you establish processes that consider users' feedback and use this to adapt the system?	2	P	EU

P: planning, R: requirement, D: design, I: implementation, T: testing, DE: deployment, O: operation



# Responsible AI Metrics Catalogue

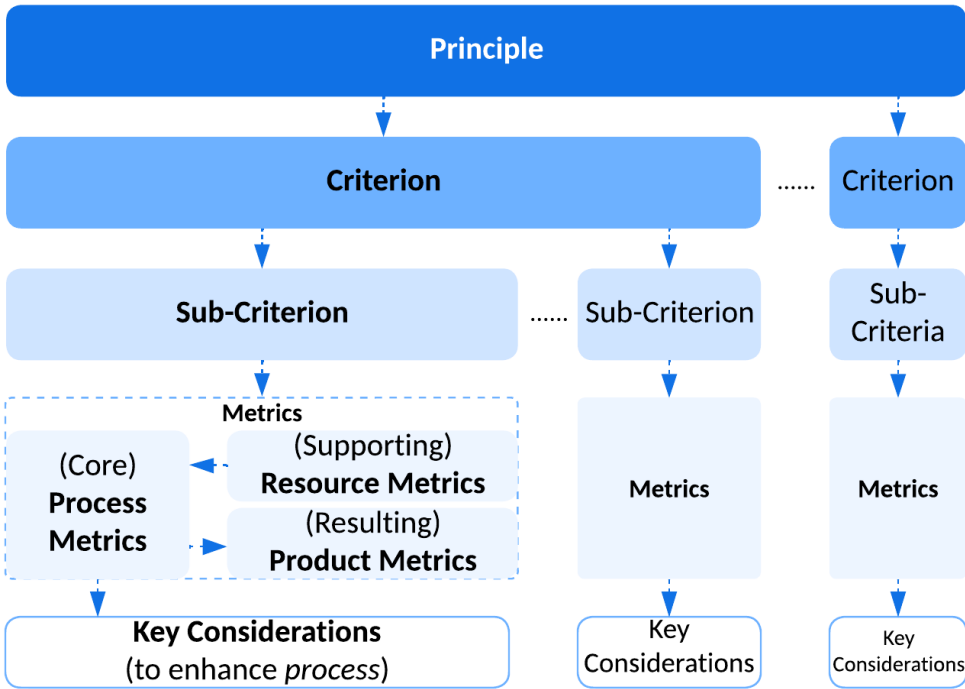


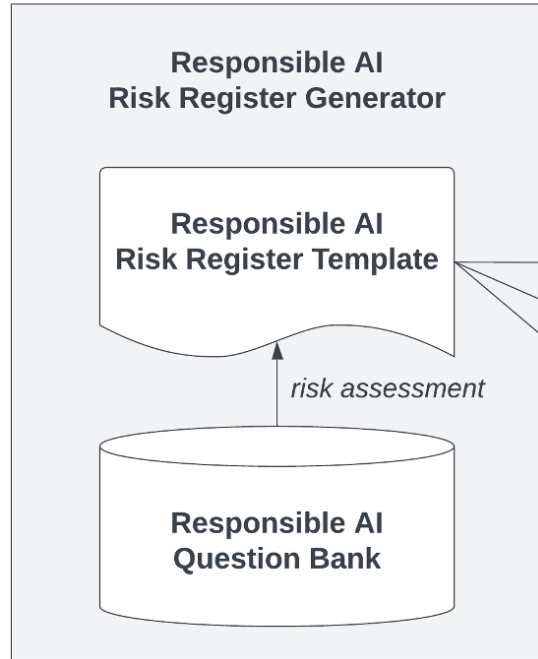
Table 2: System-Level Metrics Catalogue for AI Accountability

Criteria	Sub-Criteria	Process Metrics	Key Considerations	Resource Metrics	Product Metrics
Responsibility	RAI Oversight	Roles and Responsibilities	<ul style="list-style-type: none"> <li>Comprehensive role clarity:               <ul style="list-style-type: none"> <li>Design and development</li> <li>Deployment and operations</li> <li>Procurement and integration</li> <li>Governance and compliance</li> <li>AI as a service</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>Soft laws (e.g., best practices, guidelines standards etc)</li> <li>Hard laws (e.g., EU AI Act)</li> </ul>	<ul style="list-style-type: none"> <li>Procedure Manuals</li> <li>Contracts or agreements</li> <li>Position descriptions</li> <li>Recruitment practices</li> <li>Workforce dev strategy</li> </ul>
		RAI Governance Committee	<ul style="list-style-type: none"> <li>Multidisciplinary composition</li> <li>Strategic leadership involvement</li> </ul>		<ul style="list-style-type: none"> <li>Policy doc on Committee</li> </ul>
	Organizational AI Risk Tolerance	<ul style="list-style-type: none"> <li>Tiered risk-based categorization</li> <li>Balancing competing interests</li> <li>Holistic training content</li> <li>Targeted training for diverse roles</li> <li>Adaptive and ongoing education</li> </ul>	<ul style="list-style-type: none"> <li>Policy doc on org's risk tolerance and mitigations</li> </ul>		
	RAI Competence	<ul style="list-style-type: none"> <li>RAI Training</li> <li>RAI Capability Assessment</li> </ul>	<ul style="list-style-type: none"> <li>Training certificates</li> <li>Assessment reports</li> </ul>		
Auditability	Systematic Oversight	Data Provenance	<ul style="list-style-type: none"> <li>Detailed data record-keeping</li> <li>Data version control</li> <li>Data integrity and risk mitigation</li> <li>Legal and ethical compliance</li> </ul>	<ul style="list-style-type: none"> <li>Soft laws (e.g., auditing guidelines and frameworks etc)</li> <li>Hard laws (e.g., EU AI Act)</li> <li>AI documentation tools (e.g., datasheets, model/system cards)</li> <li>Technical tools (e.g., blockchain, knowledge graph)</li> </ul>	<ul style="list-style-type: none"> <li>Provenance records</li> <li>System features (e.g., auto-logging, version control)</li> </ul>
		Model Provenance	<ul style="list-style-type: none"> <li>Detailed model record-keeping</li> <li>Model selection and validation</li> <li>Model version control</li> </ul>		<ul style="list-style-type: none"> <li>Provenance records (and logs)</li> <li>System features (e.g., auto-logging, version control)</li> </ul>
	System Provenance and Logging	<ul style="list-style-type: none"> <li>Detailed system record-keeping</li> <li>System version control</li> <li>Decision/Trade-off</li> <li>Comprehensive operational logging</li> <li>User interaction and system response</li> <li>Incident and response</li> <li>System configuration changes</li> </ul>	<ul style="list-style-type: none"> <li>Audit reports</li> <li>Compliance certificates and licenses</li> </ul>		
	Compliance Checking	Auditing	<ul style="list-style-type: none"> <li>Composition Management</li> <li>Diversified auditing strategy</li> <li>Multi-dimensional audit techniques</li> <li>Ethical and legal compliance</li> <li>Regular audits</li> <li>Verifiable audits</li> <li>Audit-driven improvements</li> </ul>		<ul style="list-style-type: none"> <li>Incident and response doc</li> <li>System features (user feedback and report)</li> <li>System features (redundant components/functionalities)</li> </ul>
Redressability	Redress-by-Design	<ul style="list-style-type: none"> <li>Incident Reporting and Response</li> <li>Structured Incident Management</li> <li>Feedback Loop Integration</li> </ul>	<ul style="list-style-type: none"> <li>Accessibility and Visibility</li> <li>Structured Incident Management</li> <li>Feedback Loop Integration</li> </ul>	<ul style="list-style-type: none"> <li>Redundancy design case studies</li> <li>Incident management tools</li> </ul>	<ul style="list-style-type: none"> <li>Incident and response doc</li> <li>System features (user feedback and report)</li> <li>System features (redundant components/functionalities)</li> </ul>
		Built-in Redundancy	<ul style="list-style-type: none"> <li>Multi-Modal Redundancy</li> </ul>		

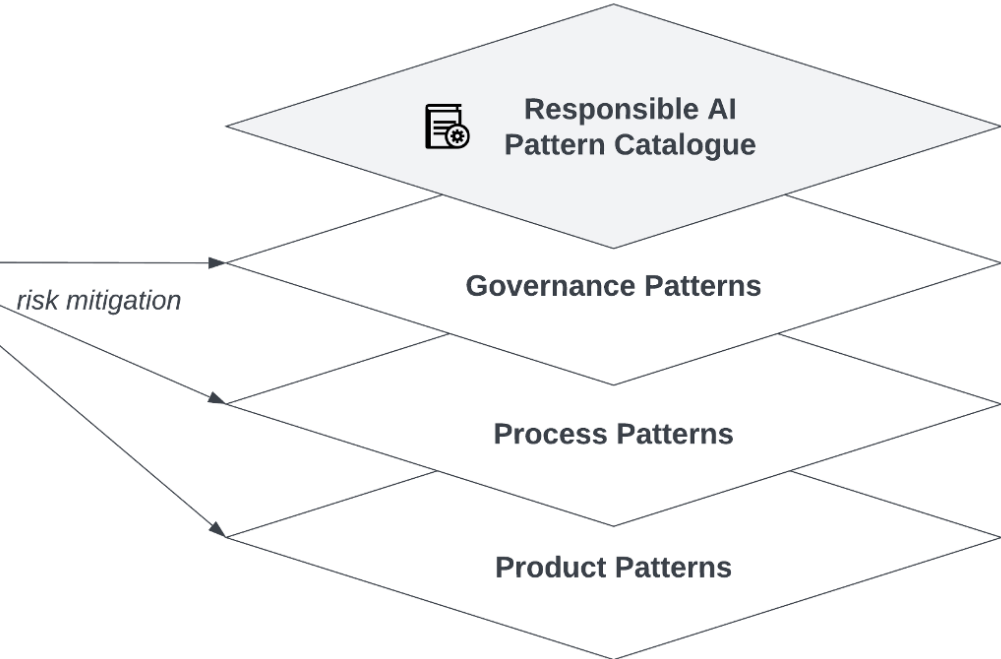


# Responsible AI Pattern Catalogue

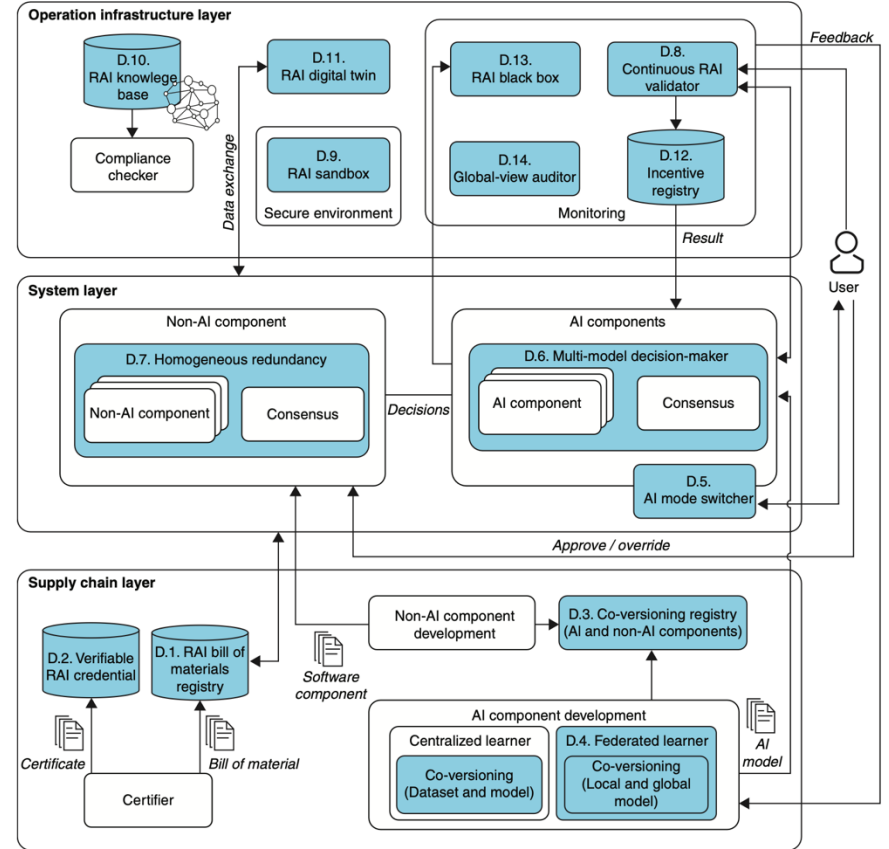
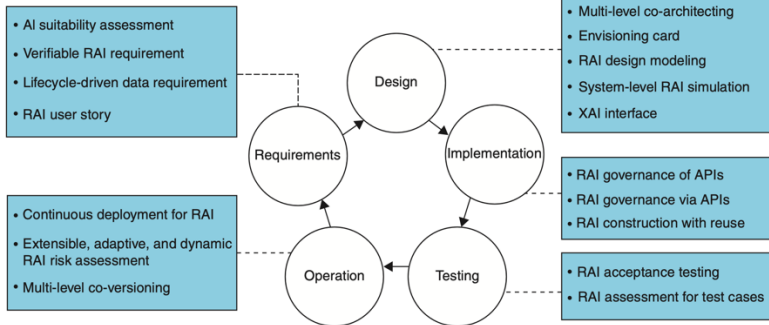
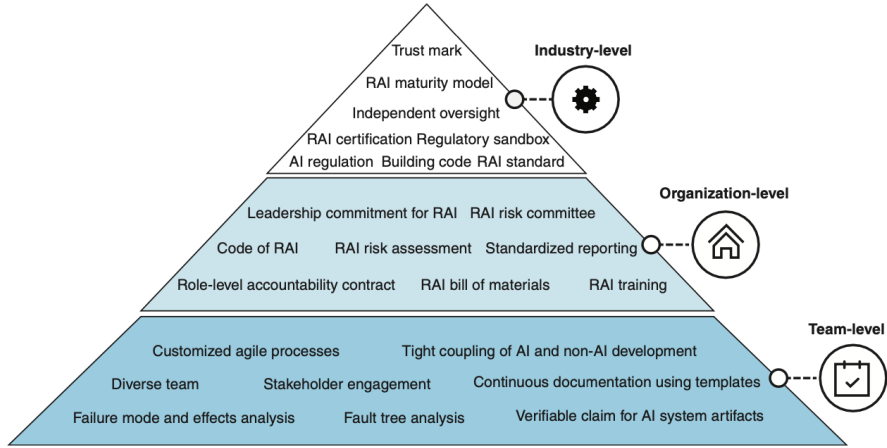
## Responsible AI Risk Assessment



## Pattern-Oriented Risk Mitigation



# Responsible AI Pattern Catalogue





## D.5. AI Mode Switcher

Adding an AI mode switcher to the AI system offers users efficient invocation and dismissal mechanisms for activating or deactivating the AI component when needed.

### Context

Human autonomy is an individual's capacity for self-determination or self-governance, which should be supported in AI systems.

### Problem

How can we enable human autonomy by allowing users to efficiently activate and deactivate the AI component when needed?

### Solution

When to use AI at decision-making points can be a major architectural design decision when designing an AI system. In Figure 6.6, adding an AI mode switcher to the AI system offers users efficient invocation and dismissal mechanisms for activating and deactivating the AI component whenever needed, thus deferring the architectural design decision to the execution time that the end user or the operator of the AI system decides. The AI mode switcher is like a *kill switch* for the AI system that could immediately shut down the AI component and thus stop its negative effects (e.g., turning off the automated driving system and disconnecting it from the internet).

### Benefits

Here are the benefits of the AI mode switcher pattern:

- **Increased trust:** An AI mode switcher gives users the choice to switch off the AI model when they do not trust the decision or recommendation provided by the AI component, thus increasing trust toward the AI system.
- **Contestability and autonomy:** The AI mode switcher enables human autonomy by allowing end users to switch off the AI component or override the decisions made by the AI component at runtime.

### Drawbacks

Here are the drawbacks of the AI mode switcher pattern:

- **Efficiency:** Efficiency and performance of the decision-making points highly depend on the quality of other non-AI components involved.
- **Suitability to (near) real-time systems:** The use of an AI mode switcher in a (near) real-time system might be problematic. The performance of the system might be affected if the end user or the operator of the AI system keeps switching the AI component on and off.

### Known Uses

Here are the known uses of the AI mode switcher pattern:

- Tesla Autopilot has multiple driver-assistance features that can be enabled or disabled during driving.<sup>28</sup> Users maintain control of the vehicles and can override the operation of these features at runtime.
- Waymo operates self-driving cars with an automated driving system that human safety drivers can override.<sup>29</sup>



# Responsible AI Chatbot



You Are Using AS



Stakeholder



Expertise

Which type of stakeholder are you?

Manager Technician Consultant Client **Board Member** Regulator Investor

Which industry sector are you from?

Health Mining Law Finance Agribusiness Cyber Security Education Defence

Infrastructure Manufacturing **R&D or Innovation** Environment

What type of AI technology is involved?

Machine learning **Language processing** Robotics Knowledge representation Computer vision

Do not know

How can I help you today?

Stakeholder

I came across this news article this morning: <https://qz.com/653084/microsofts-disastrous-tay-experiment-shows-the-hidden-dangers-of-ai> Our organisation is currently developing a chatbot for science, so I'm wondering if we might face similar AI risks?

Expertise

The final extracted key information:

Ask something

Console

AI: Now analyzing your input, please be patient.

AI: Check whether there is hyperlink.....

AI: Extracting hyperlink.....

AI: Accessing hyperlink.....

AI: Extracting Key Information from the content.....

AI: Summarizing your central ideas.....

AI: Selecting incidents from database for you.....

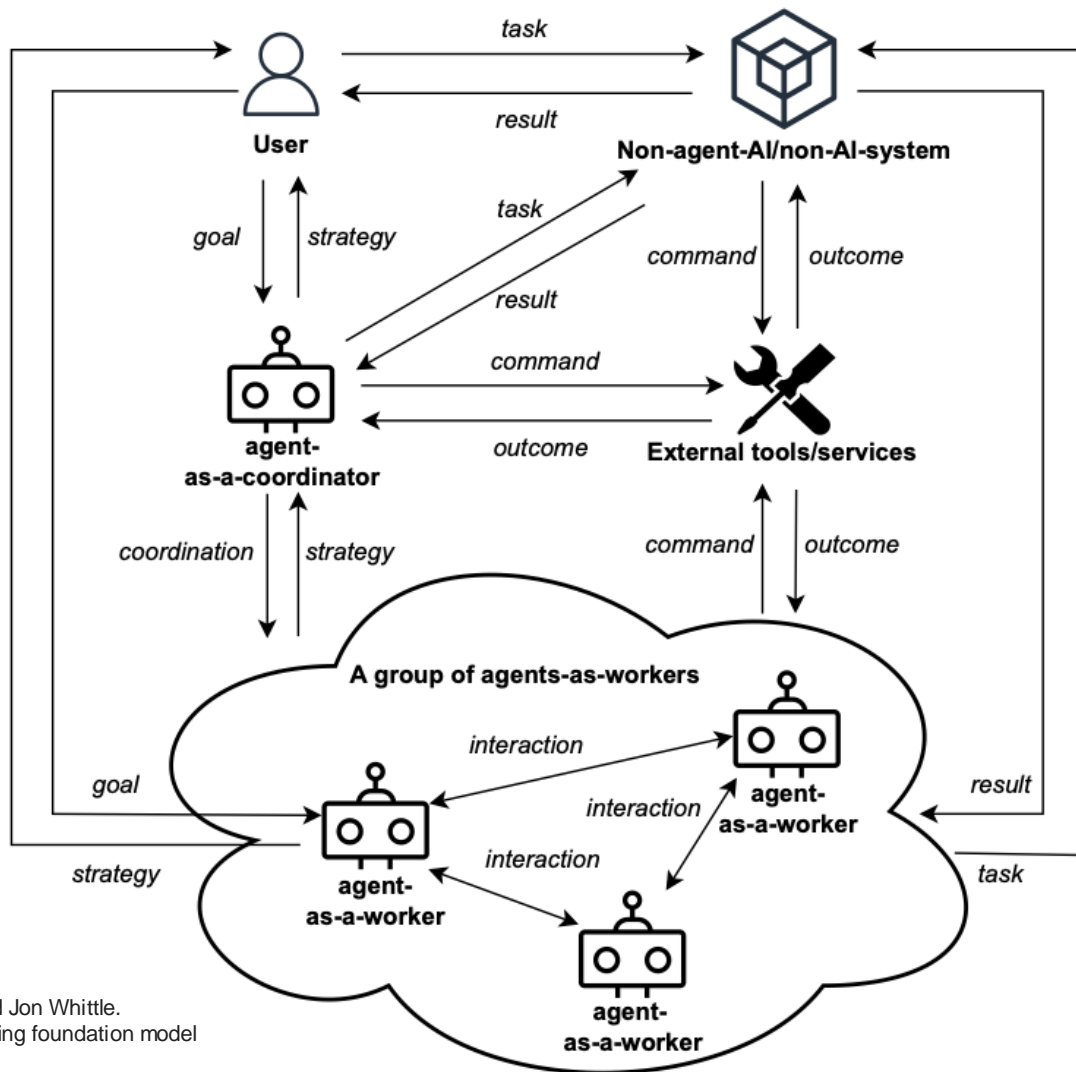
AI: Generating explanations.....



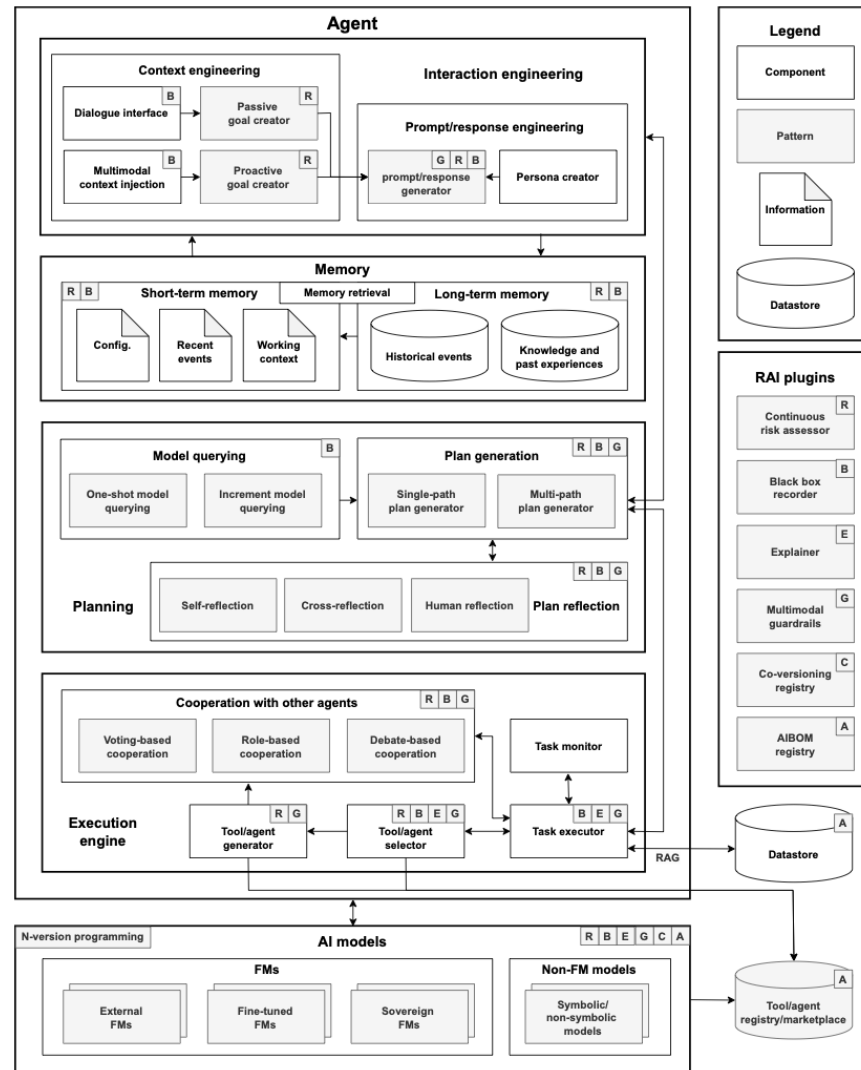
# Unique Characteristics of Agents

- Complex Architecture
- Autonomous Operation
- Non-Deterministic Behaviour
- Continuous Evolution

# Architecture of an agent-based ecosystem



# Reference architecture

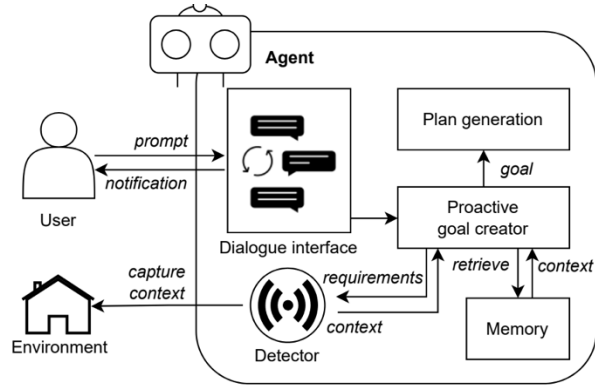


Qinghua Lu, Liming Zhu, Xiwei Xu, Zhenchang Xing, Stefan Harrer, and Jon Whittle. "Towards responsible generative ai: A reference architecture for designing foundation model based agents." ICSA'24. <https://arxiv.org/pdf/2311.13148>



# Agent Design Pattern Catalogue

## Pattern: Proactive Goal Creator



**Summary:** Proactive goal creator anticipates users' goals by understanding human interactions and capturing the context via relevant tools.

**Context:** Users explain the goals that the agent is expected to achieve in the prompt.

**Problem:** The context information collected via solely a dialogue interface may be limited, and result in inaccurate responses to users' goals.

### Forces:

- **Underspecification.** i) Users may not be able to provide thorough context information and specify precise goals to agents. ii) Agents can only retrieve limited information from the memory.
- **Accessibility.** Users with specified disabilities may not be able to directly interoperate with the agent via *passive goal creator*.

**Solution:** Fig. 4 illustrates a simple graphical representation of *proactive goal creator*. In addition to the prompts received from dialogue interface, and relevant context retrieved from memory, the *proactive goal creator* can anticipate users' goals by sending requirements to detectors, which will then capture and return the user's surroundings for further analysis and comprehension to generate the goals, for instance, identifying the user's gestures through cameras, recognising application UI layout via screenshots, etc. The *proactive goal creator* should notify users about context capturing and other relevant issues with a low false positive rate, to avoid unnecessary interruptions. In addition, the captured environment information can be stored in the agent's memory (or knowledge base) to establish "world models" [22, 23] to improve its ability to comprehend the real world.

### Consequences:

#### Benefits:

- **Interactivity.** An agent can interact with users or other agents by anticipating their decisions proactively with captured multimodal context information.
- **Goal-seeking.** The multimodal input can provide more detailed information for the agent to understand users' goals, and increase the accuracy and completeness of goal achievement.
- **Accessibility.** Additional tools can help capture the sentiments and other context information from disabled users, ensuring accessibility and broadening the human values of foundation model-based agents.

#### Drawbacks:

- **Overhead.** i) *Proactive goal creator* is enabled by the multimodal context information captured by relevant tools, which may increase the cost of the agent. ii) Limited context information may increase the communication overhead between users and agents.

### Known uses:

- **GestureGPT** [24]. GestureGPT can decipher users' hand gesture descriptions and hence comprehend users' intents.
- Zhao et al. [25] proposed a programming screencast analysis tool that can extract the coding steps and code snippets.
- **ProAgent** [26]. ProAgent can observe the behaviours of other teammate agents, deduce their intentions, and adjust the planning accordingly.

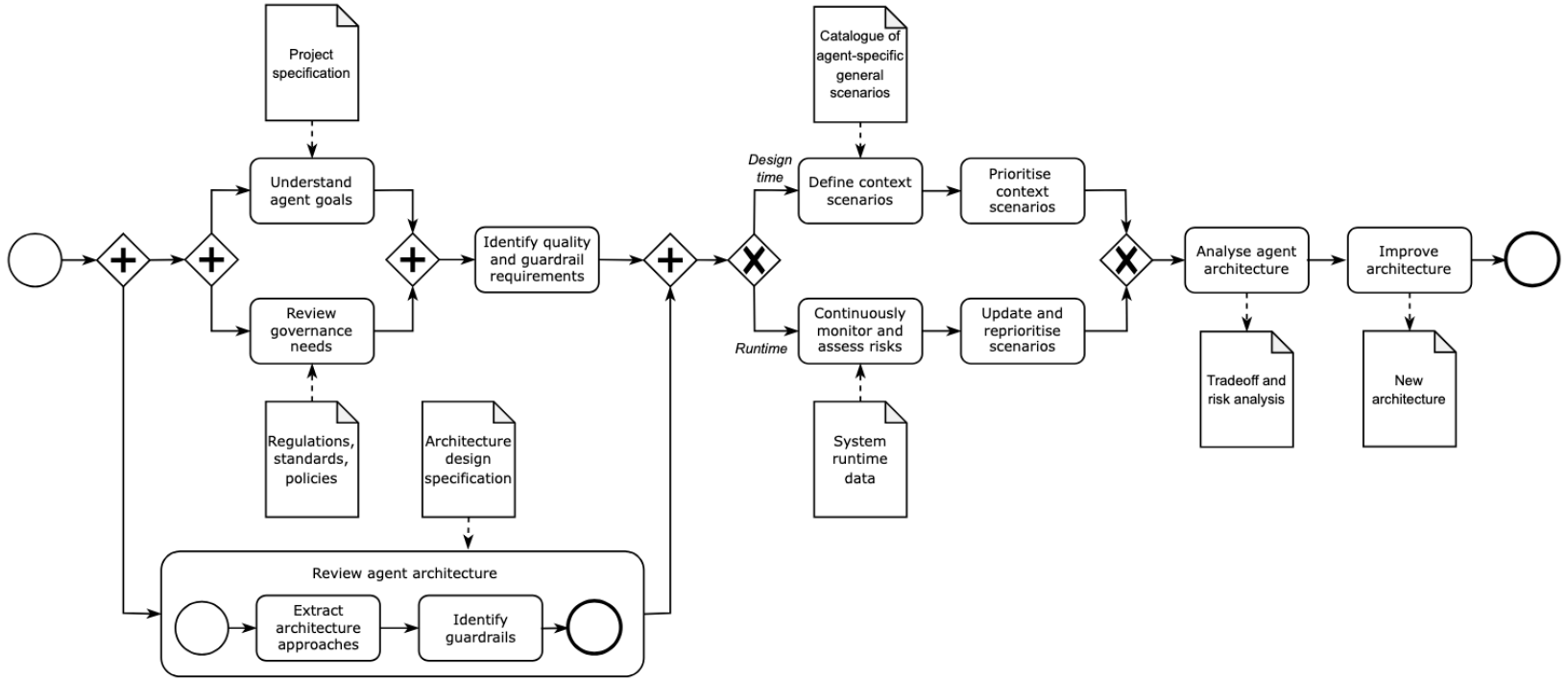
### Related patterns:

- **Passive goal creator.** *Proactive goal creator* can be regarded an alternative of *passive goal creator* enabling multimodal context injection.
- **Prompt/response optimiser.**
- *Proactive goal creator* can first handle users' inputs and transfer the goals and relevant context information to *prompt/response optimiser* for prompt refinement.





# AgentArcEval: An Architecture Evaluation Method for Foundation Model based Agents



Qinghua Lu, Dehai Zhao, Yue Liu, Hao Zhang, Liming Zhu, Xiwei Xu, Angela Shi, and Tristan Tan. "AgentArcEval: An Architecture Evaluation Method for Foundation Model based Agents." (2024).  
[https://www.researchgate.net/publication/385660422\\_Evaluating\\_the\\_architecture\\_of\\_large\\_language\\_model\\_based\\_agents](https://www.researchgate.net/publication/385660422_Evaluating_the_architecture_of_large_language_model_based_agents)

# Catalogue of Agent-Specific General Scenarios

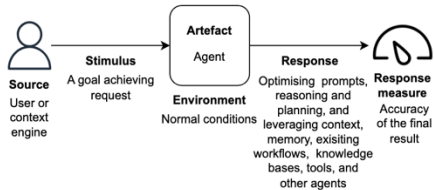


Fig. 2: Accuracy general scenario.

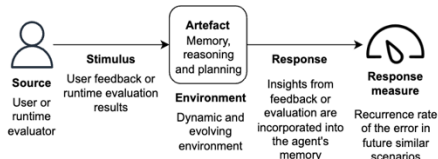


Fig. 3: Adaptability general scenario.

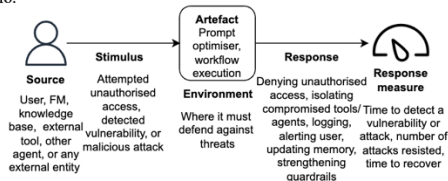


Fig. 6: Security general scenario.

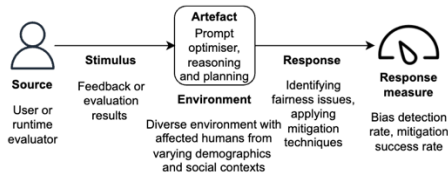


Fig. 7: Fairness general scenario.

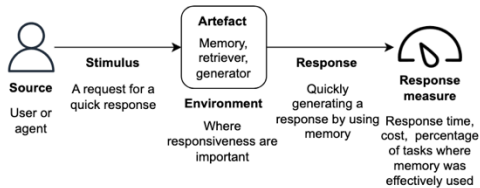


Fig. 4: Efficiency general scenario.

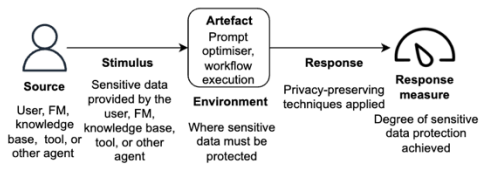


Fig. 5: Privacy general scenario.

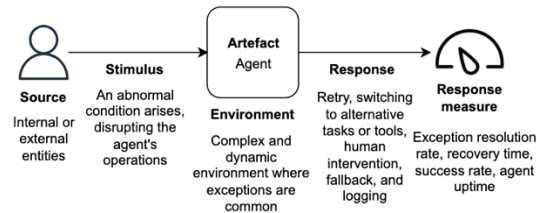


Fig. 8: Reliability general scenario.

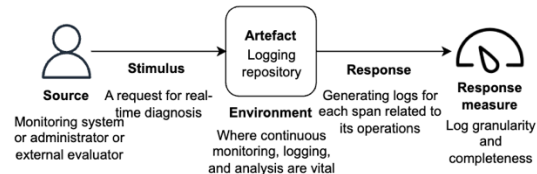


Fig. 9: Observability general scenario.

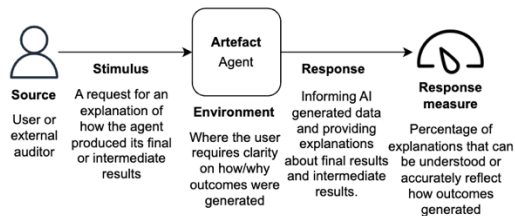
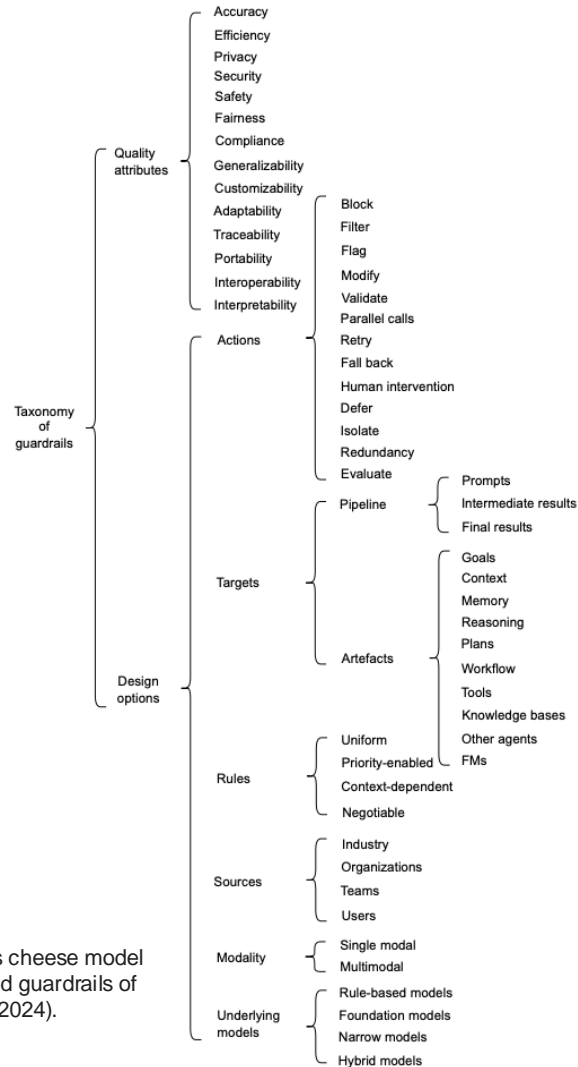


Fig. 10: Transparency general scenario.

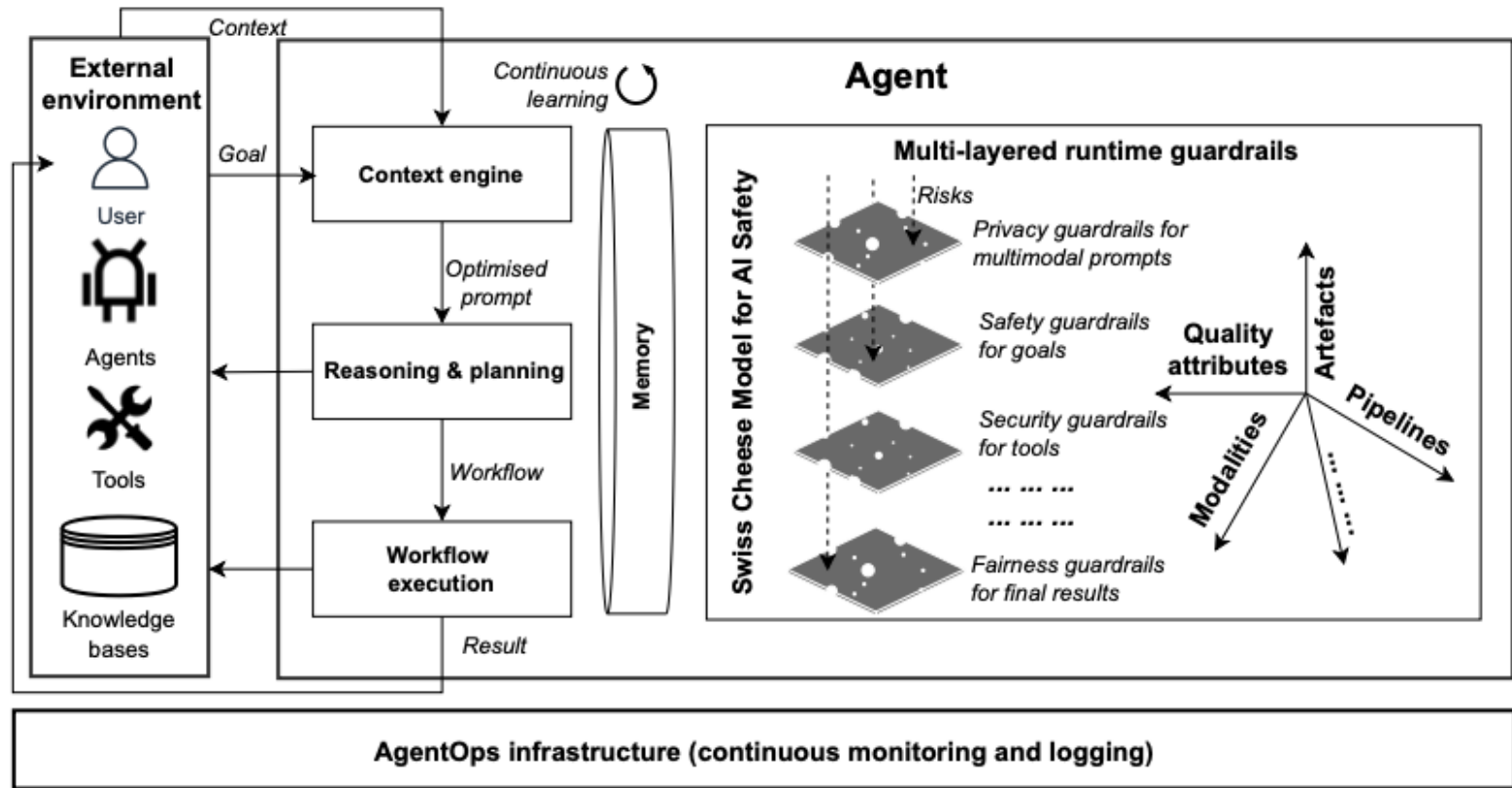


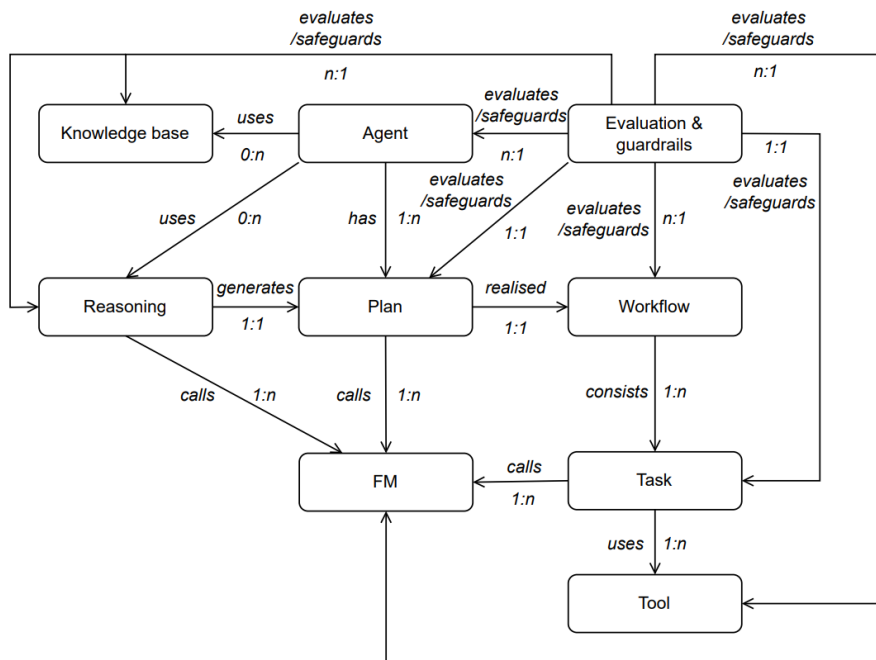
# Design Taxonomy of Runtime Guardrails for LLM Agents



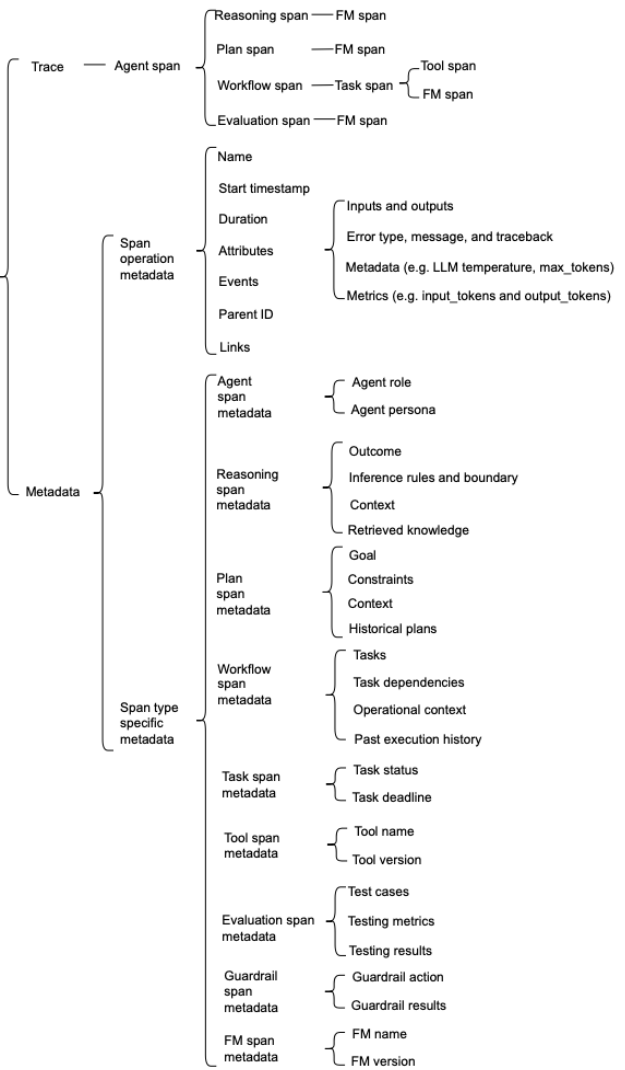
Md Shamsujjoha, Qinghua Lu, Dehai Zhao, and Liming Zhu. "Swiss cheese model for ai safety: A taxonomy and reference architecture for multi-layered guardrails of foundation model-based agents." *arXiv preprint arXiv:2408.02205* (2024).  
<https://arxiv.org/abs/2408.02205>

# Swiss Cheese Model for AI Safety – Multi-Layered Guardrails for LLM Agents

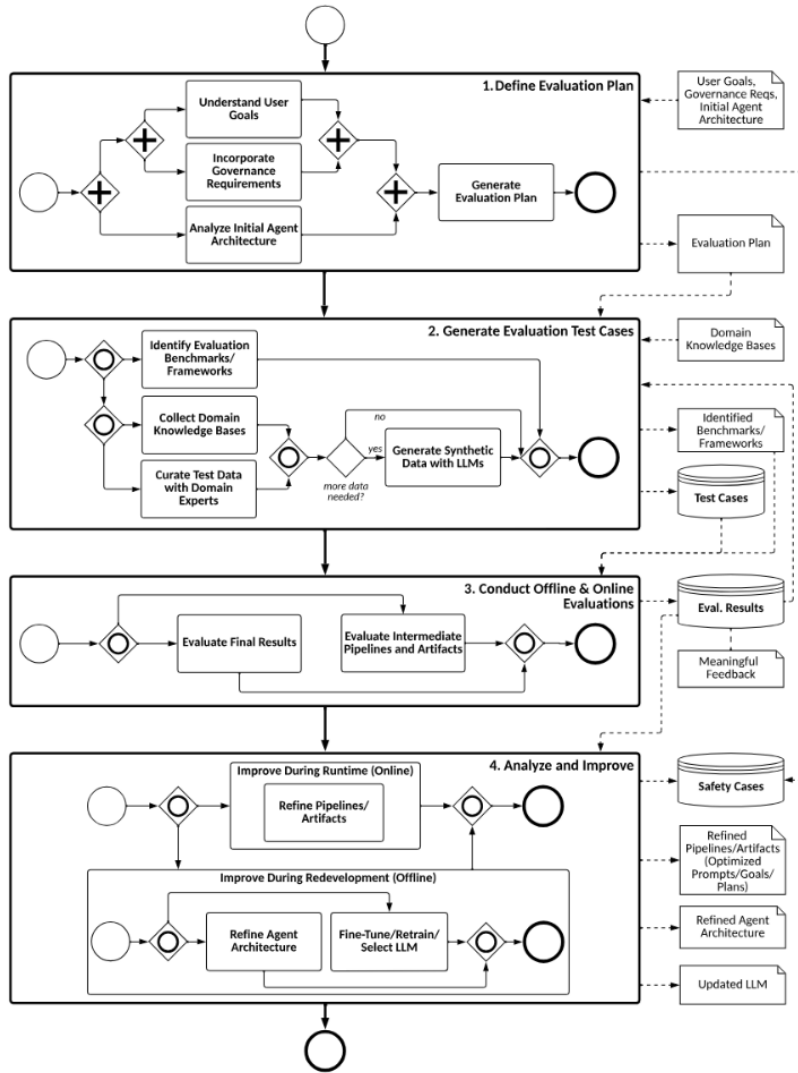




Taxonomy of AgentOps

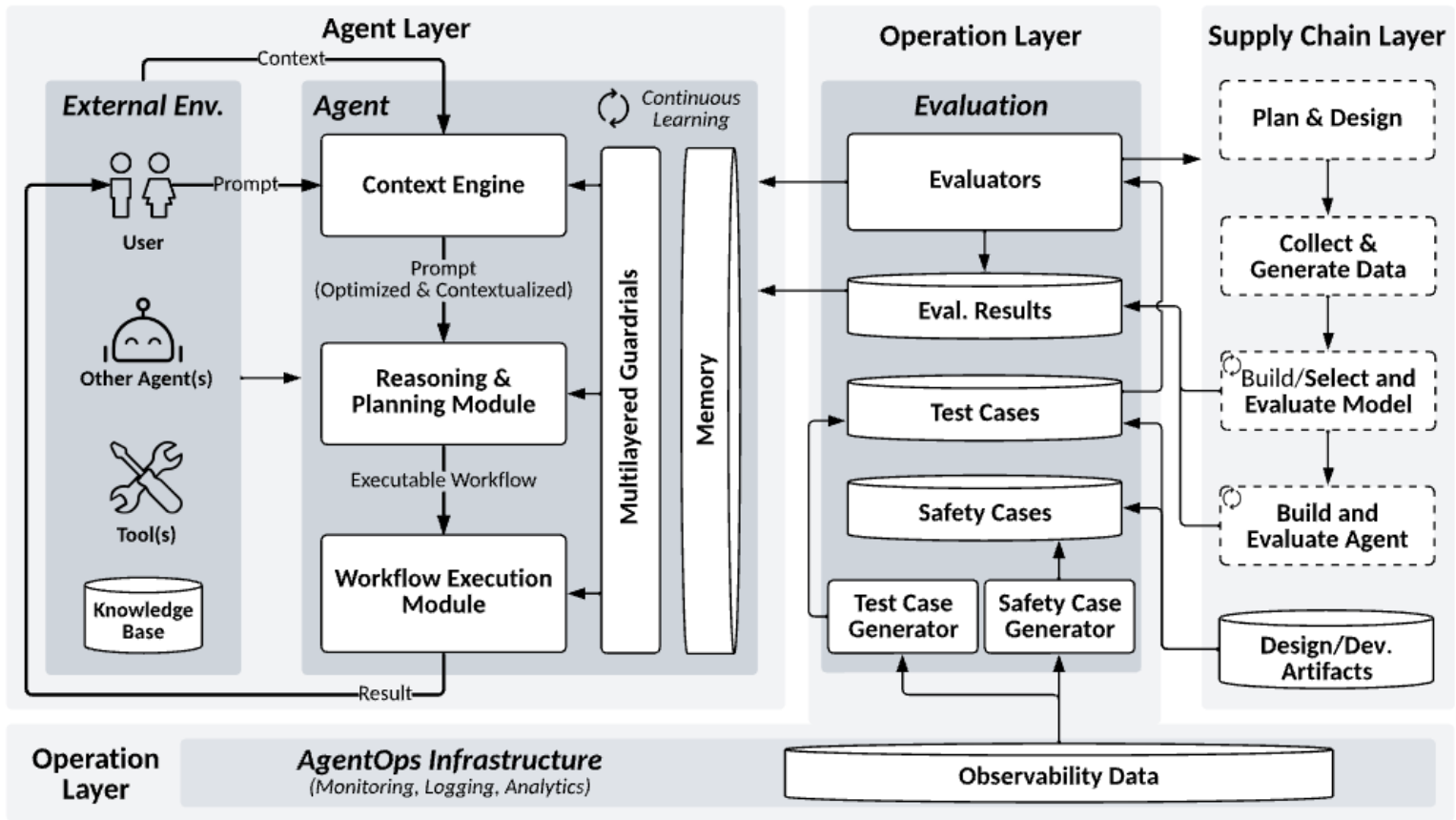


# Process Model for LLM Agent Evaluation





# Evaluation-Driven Design for LLM Agents





# Science Digital: Agent Platform

Navigation: Back | 1 SPL Crafting | 2 SPL Testing | Next

**SPL Form** | Save Changes | Linting → | Workflow Inspection | Sapper Copilot

**Persona** The Agentware you need the AI native service to play

Description: A smart shopping agent with product browsing capabilities.

**Audience** The target audience you want the Agentware to serve

Description: Online shoppers looking for personal recommendations.

**Terminology** For some specific nouns

Term: Product browsing

Term: Shopping agent

Model: gpt-4-1106-preview

Temperature (Coming soon): 1000

Max length (Coming soon): 8000

Knowledge files (Coming soon) Add File

**Sapper Copilot**

Hi, I am Sapper Copilot. I can help you generate, analyze, and revise SPL Form.

I am looking to create a shopping agent with a function agent: product browsing.

Working on it...

Initializing agent...Creating SPL form...

Navigation: Back | ✓ SPL Crafting | 2 SPL Testing

AL Chain | SPL Prompt | Website | WeChat | Robot

Flowchart:

```
graph TD; A[Wait for user input message.] --> B[Customer support]; A --> C[Product browsing]; B --> D[Output the customer service answer.]; C --> E[Output the filtered products.];
```





# Responsible/Safe AIware Engineering at system-level and across supply chain

Thank you.

Qinghua Lu

Responsible AI Science Team Leader

[qinghua.lu@data61.csiro.au](mailto:qinghua.lu@data61.csiro.au)

<https://research.csiro.au/ss/team/se4ai/responsible-ai-engineering/>

