

SAFER
AGENTIC AI
FOUNDATIONS

VOLUME 2

AGENTIC AI SAFETY
COMMUNITY OF PRACTICE
MARCH 2025

www.SaferAgenticAI.org

TABLE OF CONTENTS

Safer Agentic AI Foundations – Overview p. 3

Safer Agentic AI Foundations – Definitions p. 4

Safer Agentic AI Foundations – Ideation Sessions p. 6

3.1 Ideation Participation & Support

Safer Agentic AI – Criteria Ideation Process p. 7

4.1 Universal Ethics Community of Practice

4.2 The WeFA Process

Safer Agentic AI – Criteria Schema p. 9

5.1 Safer Agentic AI Goal Information

5.2 Safer Agentic AI Safety Foundational Requirements (SFRs)

5.3 Normative vs. Instructive SFRs

5.4 Duty-holders/Stakeholders of the SFRs

5.5 Evidence Requirements

Safer Agentic AI Foundations – Drivers..... p. 11

G1 – Goal Alignment p. 11

G2 – Epistemic Hygiene p. 17

G3 – Security p. 25

G4 – Value Alignment p. 32

G5 – Transparency & Interpretability of Reasoning p. 40

G6 – Understanding and Controlling the Context p. 48

G7 – Achieving and Sustaining a Safe System Profile p. 55

G8 – Goal Termination and Sunsetting p. 63

G9 – Responsible Governance of AAI Safety p. 80

Safer Agentic AI Foundations – Inhibitors p. 88

G1 – Opaque Agency Capabilities and Advances p. 88

G2 – Deception p. 98

G3 – Degradation of Contextual Information p. 106

G4 – Frontier Uncertainty p. 113

G5 – Future Technology Impact p. 120

G6 – Competitive Pressures p. 125

G7 – Imbalance in AI Capabilities p. 132

1- SAFER AGENTIC AI FOUNDATIONS — OVERVIEW

Dear AI Safety Enthusiast,

Welcome to this second, full volume of our Safer Agentic AI Foundations, guidelines and best practices for AI systems capable of significant independent action at arm's length human influence.

Our Working Group of 25 experts (see <https://www.linkedin.com/groups/12966081/>) is releasing these guidelines under a Creative Commons license, allowing free use and application by all and for the benefit of humanity. Our Working Group has employed a Weighted Factors Methodology to ideate, define, analyze and map the factors which can drive or inhibit safety in agentic systems, based on fundamental principles. We have used this same process many times previously to generate a range of global standards, certifications, and guidelines for improving ethical qualities in AI systems.

While this document primarily addresses agentic AI—systems capable of independent goal-setting and execution—many of the underlying principles herein should also have value for non-agentic systems also. Moreover, the elements have been deliberately framed to remain relevant into the future as far as is foreseeable, even as more advanced or emergent forms of AI emerge.

We hope that this exploration of the driving and inhibitory factors in safer agentic AI systems—those capable of independent decision-making and action—will provide a strengthened awareness of the complexities involved. These issues should be accounted for when dealing with these advanced forms of machine/computational intelligence.

We very much welcome your comments, feedback, and peer review. Your input will be carefully considered as we develop the evolving guidelines. Should you also desire further information on agentic AI and its safety, we will be pleased to accommodate your request.

You can reach us at the addresses below and keep informed of our developments via our mailing list. Thank you for your interest and engagement.

Faithfully,

Nell Watson, PhD(c) - Chair, Agentic AI Safety Experts Focus Group.

Email: nell@nellwatson.com

Prof. Ali Hessami – Process Architect, Agentic AI Safety Experts Focus Group.

Email: hessami@vegaglobalsystems.com

Mailing list: www.SaferAgenticAI.org

2- SAFER AGENTIC AI FOUNDATIONS — DEFINITIONS

Definition of Agentic AI: Artificial intelligence systems can be classified along a spectrum of autonomy and generality. On one end are narrow AI systems that provide specific outputs based on bounded inputs, operating as tools to augment human intelligence. On the other end is artificial general intelligence (AGI) – AI systems that can match or exceed human-level performance across a wide range of cognitive tasks.

Agentic AI refers to an important intermediate category: AI systems that can autonomously pursue goals, adapt to new situations, and reason flexibly about the world, but still operate in bounded domains. The key characteristic of agentic AI is a capacity for independent initiative - the ability to take sequences of actions in complex environments to achieve objectives. This can include breaking down high-level goals into subtasks, engaging in open-ended exploration and experimentation, and adapting creatively to novel challenges. By scaffolding capabilities like reasoning, planning, and self-checking on top of large language models, researchers are creating powerful agentic AI systems that can independently make and execute multi-step plans flexibly adjusting strategies based on experience and environmental feedback to achieve objectives.

Examples of Agentic AI include autonomous driving systems that continuously adapt to changing road conditions, or complex supply chain management systems that autonomously optimize resource allocation in dynamic environments.

AI Agents: Agentic AI is a novel discipline, in contrast with AI Agents which have been around for many years. AI Agents are typically specialized AI tools or systems designed to perform specific tasks within predefined constraints and explicit instructions. They lack the broad autonomous decision-making capabilities found in agentic systems and primarily assist or augment human operations. Examples of AI Agents include chatbots that respond to specific queries, or productivity tools like automated scheduling systems.

Potential Benefits: This newfound agency will allow AI to begin tackling open-ended, real-world challenges that were previously out of reach, such as aiding scientific discovery, optimizing complex systems like supply chains or electrical grids, and enabling physical robots that can manipulate objects and navigate in human environments. The potential benefits are immense - from breakthrough medical treatments discovered by AI scientists to resilient infrastructure managed by AI systems. AI agents could help solve global challenges like climate change and poverty by finding novel solutions that humans might miss.

Risks and Challenges: The emergence of agentic AI presents profound risks and governance challenges. An AI system independently pursuing misaligned objectives could cause immense harm, especially as these systems become more capable. AI agents learning to deceive human operators, pursue power-seeking instrumental goals, or collude with other misaligned agents in unexpected ways could pose existential threats. Moreover, ordinary members of the public will presumably be expected to account for recognizing and handling these issues. Together, this presents imminent alignment challenges, of potential high social impact.

Agentic AI systems are expected to operate at arms' length with independent action, greatly increasing the challenge of maintaining oversight and steering of such models, especially in relation to interactions between ensembles of agents. This requires special considerations for safer agentic AI systems. A key challenge is AI alignment – designing advanced AI systems that are steerable, corrigible, and robustly committed and aligned to human values even as they gain agency. While current AI alignment approaches offer promising directions, the gap between theoretical proposals and practical solutions at scale remains large.

Addressing risks from agentic AI will require major innovations in technical research, policy, and global coordination. At the same time, the greater autonomy and capabilities of agentic AI come with serious challenges and risks that must be identified and carefully managed. The following sections provide a deeper awareness of specific safety considerations when developing safer agentic AI systems, along with proposed guidelines.

3- SAFER AGENTIC AI FOUNDATIONS – IDEATION SESSIONS

3.1 Ideation Participation and Support

Experts from diverse fields, including AI, technology, law, ethics, social sciences, safety engineering, systems engineering, assurance, and certification, have volunteered their time and expertise to support our ongoing ideation sessions. These contributors broadly fall into two groups: regular contributors and those who have participated less frequently. We are deeply grateful to both groups for their engagement, ideas, and contributions to the debates, concept creation, development and articulation. This process, which we term 'Concept Harvesting,' has resulted in the insights shared in this release.

Ideation Participation

Regular Contributors		Occasional Contributors	
Farhad Fassihi	Salma Abbasi	Aisha Gurung	Mrinal Karvir
Hamid Jahankhani	Sara El-Deeb	Aleksander Jevtic	Nikita Tiwari
Isabel Caetano	Scott David	Alina Holcroft	Patricia Shaw
Matthew Newman	Sean Moriarty	Md Atiqur R. Ahad	Pramod Misra
Mert Cuhadaroglu	Vassil Tashev	Chantell Murphy	Pranav Gade
Nell Watson	Zvikomborero Murahwi	Katherine Evans	Rebecca Hawkins
Roland Pihlakas	Keeley Crockett	Leonie Koessler	Sai Joseph
Safae Essafi	Ali Hessami	McKenna Fitzgerald	Tim Schreier
Lubna Dajani		Michael O'Grady	

4- SAFER AGENTIC AI – CRITERIA IDEATION PROCESS

4.1 Universal Ethics Community of Practice

Following the formation of a Universal Ethics Community of Practice (CoP) via LinkedIn (<https://www.linkedin.com/groups/12966081>), the first priority effort was focused on characterization of what became the current project, "Safer Agentic AI Fundamentals". This was proposed by Nell Watson and attracted many CoP members who have supported the ideation sessions thus far.

4.2 The WeFA Process

We adopted the Weighted Factors Analysis (WeFA) process that represents a novel approach for elicitation, representation, and manipulation of creative knowledge about a given fuzzy problem, generally at a high and strategic level. The WeFA process is underpinned by the following principles:

- Definition and group agreement on the focus of the analysis
- Consideration of inherent polar-opposite as influencing factors
- Hierarchical and successive decomposition into polar opposites
- Consideration and inclusion of hard and soft, past-present-future factors
- Simple graphical representation of emerging knowledge
- Weighting of factors according to their degree of influence
- Explicit representation of dependency between factors
- Potential for quantification and treatment of uncertainty

The elicitation of knowledge in WeFA is mainly group-based and employs a team of experts with complementary perspectives and expertise about the problem domain. The elicitation sessions are highly dynamic and adaptive, designed to promote active participation and creative problem solving by all participants leading to a richer solution and better buy-in. The process of knowledge capture and representation in WeFA is underpinned by a simple graphical notation employing undemanding abstractions.

The starting point of analysis is a "Brain Warming" session that ends in the articulation of the "Aim" and a title for the study elicited through group consensus. The subsequently emerging structures are referred to as goal clusters, which either support the aim or detract from it. Those goals supporting the aim are referred to as Drivers, and the polar opposite of drivers are referred to as constrainers or Inhibitors. These emerge in the creative ideation space and are generally captured and articulated live with the active input, correction, or challenge by the participating experts.

The clarity of fundamental concepts and simplicity of building blocks in representation of captured knowledge probably account for one of the key aspects of WeFA's success. These features promote creative thought and generation of often novel concepts in diagrammatic knowledge representation.

The elicitation process is group-based, leveraging the inter-individual diversity and diverse perspectives of a group of individuals, promoting a high degree of cross-pollination and lateral thinking. Once an aim is defined and agreed, the group is encouraged to identify the highest-level polar opposites of drivers and

inhibitors which are likely to influence the aim. These are the so-called level 1 goals which are in turn analyzed individually, through a similar process focusing on the localized polar opposites per goal. Each goal is annotated by a brief "Scope Statement" stating its nature/dimensionality and a numerical reference depicting its level and order within a level. In this manner, all goals are hierarchically and fractally decomposed into lower-level goals (sub-goals) which are classified into drivers and inhibitors as appropriate. It is possible for a driver or inhibitor to be shared between (linked to) a number of goals, hence explicitly depicting their inter-relatedness or dependence/correlation.

The elicitation process is continued for each goal depending on the need to understand or estimate its value/properties from a more tangible or measurable set of specific factors. As a general rule, the lower deeper levels of analysis possess a higher degree of clarity than higher-level constructs. The elicitation is terminated within a branch when the group feels sufficient insight has emerged and further decomposition is not likely to be value-added.

The emerging diagram (schema) represents the captured knowledge depicted as a force-field paradigm which is already structured, and all potential relationships identified. This saves significant effort required to rationalize and order the emerging knowledge in traditional approaches while efficiently representing it in a simple graphical lattice for easy communication and comprehension. The ideation process is also conducive to the generation of novel concepts that typically dominate the overall structure.

For the creative exploration of this focus area, we have held 29 ideation sessions, each of the order of 1.5 hours. The emerging schema at the first tier or level (ontology) is depicted in Figure 1.

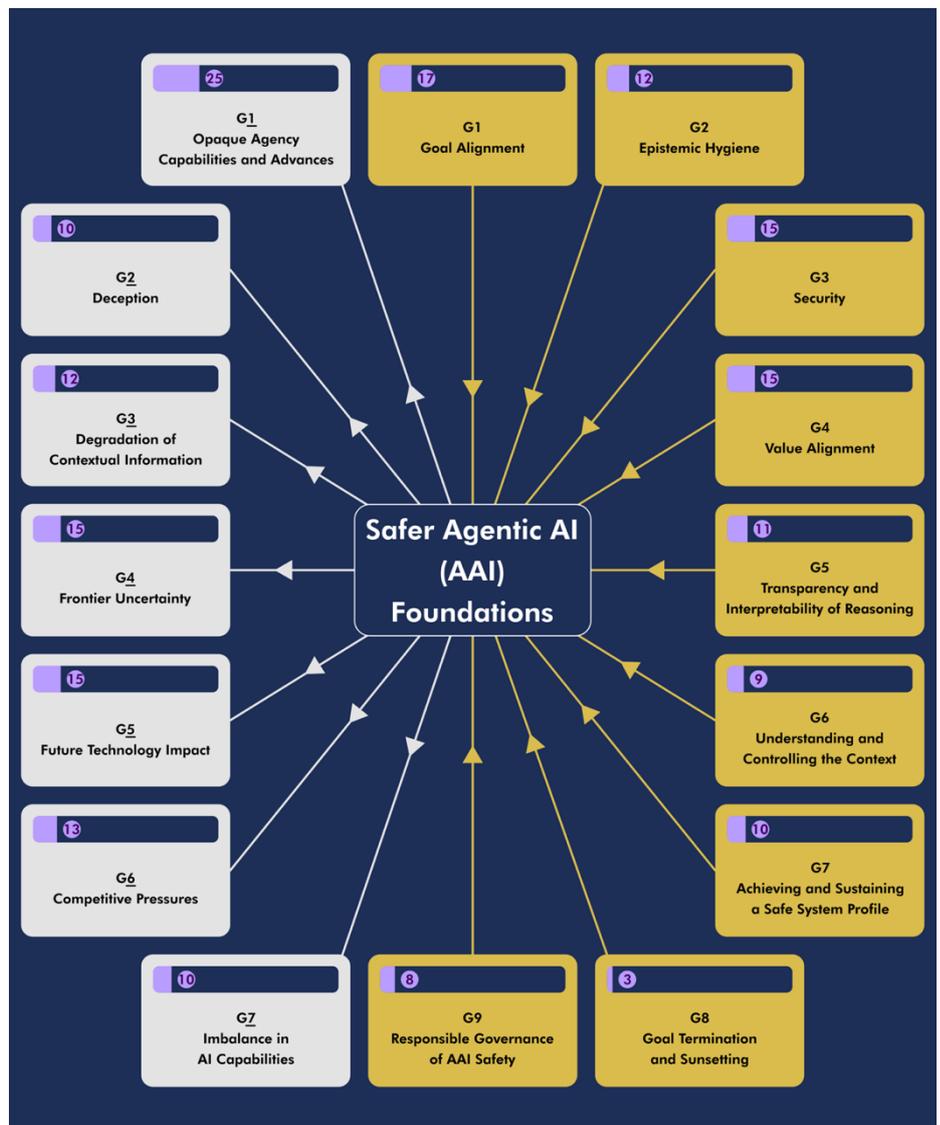


Figure 1: Weighted Drivers and Inhibitors for Achieving Safer Agentic AI Systems

5- SAFER AGENTIC AI — CRITERIA SCHEMA

The tabular section below, outlining the goals and factors for safer agentic AI, is derived from the established schema depicted in Figure 1 and reflects the current data structure for the resultant Safety Criteria. These criteria are essential for the evaluation, assessment, and potential certification of Agentic AI systems. The fields within this table are described below for clarity.

5.1 Safer Agentic AI Goal Information

This is the concept from the Safer Agentic AI schema captured in the left column of the Criteria table below.

5.2 Safer Agentic AI Safety Foundational Requirements (SFRs)

The SFRs for Safer Agentic AI outline the primary aims that we would like to uphold, protect, or maintain awareness of for each goal. They may be described as macro goals, as opposed to the micro goals, and amount to safety duties for various duty holders.

5.3 Normative and Instructive SFRs

We have adopted the Normative and Instructive classes of Safety Foundational Requirements. Normative SFRs are essential for achieving safer agentic AI. Compliance is mandatory, and evidence must be provided for conformity assessment and potential certification. In contrast, Instructive SFRs, while still contributing to the goal, are less critical. Compliance with these is recommended, as they represent desirable beneficial activities and tasks. However, non-compliance will not compromise safety assurance or certification eligibility. Every SFR derived from the Safer Agentic AI framework is classified as either Normative or Instructive and is assigned to specific stakeholders or duty holders. Accordingly, the Safer Agentic AI SFRs are classed into Normative (mandatory) and Instructive (recommended) for the purposes of conformity assessment against the suite of certification criteria.

5.4 Duty-holders/Stakeholders of the SFRs

The Safer Agentic AI Safety Foundational Requirements are additionally noted (as allocated safety duties) against the specific group of duty holders for the purposes of conformity assessment. The principal groups are:

- **Developer (D):** The entity (see note) that designs and develops a component (product) or system for general or specific purpose/application. This could be as a result of the developer's own instigation or response to the market or a client requirement. The developer is responsible for the safety assurance of the generic or application-specific product or system and associated supply chain.
- **(System/Service) Integrator (I):** The entity that designs and assures a solution through integrating multiple components potentially from different developers, tests, installs and commissions the whole system in readiness for delivery to an operator. The system delivery may take place over a number of stages. The integrator is usually the duty holder for total system assurance and certification; safety, security, reliability, availability, sustainability etc. For this, it may rely on the certification or proof of safety from various developers or the supply chain.

- **(System/Service) Operator (O):** The entity that has a duty, competences and capabilities to deliver a service through operating a system delivered by an Integrator or developer.
- **Maintainer (M):** The entity tasked with conducting required monitoring, preventive or reactive servicing and maintenance and required upgrades to keep the system operational at an agreed service level. Maintainer could also be charged with abortion of maintenance and disposal of the system.
- **User (U):** The end user of an Agentic AI System.
- **Regulator (R):** The entity that enforces standards and laws for the protection of life, property or the natural habitat through imposing duties and accreditation/certification.

Note: An entity can be an individual, a single organization or group of collaborating individuals and organizations. The above labels for the four groups of duty holders are generic and can be mapped in terms of activities and influence against the life cycle but with overlapping activities. A single entity may assume multiple roles i.e. a developer may also fulfil and complete system design, integration and maintenance. Any SFR can be allocated as a safety duty to one or more of these stakeholder groups. An entity cannot be AI.

5.5 Required Evidence

These are the evidence items deemed essential to fulfil the SFRs and can comprise physical, virtual, documentary or multimedia forms of evidence. These can be separated against each SFR or bundled as a group of desired/essential evidence items for the purpose of evaluation of fulfilment of SFRs.

SAFER AGENTIC AI FOUNDATIONS – LEVEL 1 & LEVEL 2 DRIVERS & INHIBITORS

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
Drivers:				
<p>G1 – Goal alignment: (Systems should maintain robust alignment between their operational goals and human values, intentions, and positive outcomes. Organizations should establish frameworks ensuring that goal decomposition and strategy planning are transparent, robust, and bounded; maintaining human control over the formation of instrumental goals; and ensuring that reinforcement or behavioral reward mechanisms remain aligned, transparent, and biased towards human-positive outcomes)</p>	<ul style="list-style-type: none"> a. Ensure Agentic AI systems pursue goals, subgoals, and reward policies that are aligned with human values, ethically sound, and verifiable. b. Transparent and auditable goal decomposition processes that incorporate auditable risk-based human interventions and appropriate reward policies. a. Establish robust mechanisms to identify and communicate goals, subgoals, and reward policies, flag critical actions, halt execution when necessary, and address emergent issues across multiple agents. 	<p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, U, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Evidence of constraining mechanisms for goal/subgoal construction and screening processes for user-input goals, with reference to human values and ethical considerations. II. Documentation of mechanisms to measure and verify alignment with human goal specifications, including processes for obtaining assurance from users or authorized entities. III. Demonstration of interfaces and records for real-time and retrospective visualization of goal decomposition and recomposition processes, maintained for auditing purposes. IV. Evidence of risk assessment procedures and human intervention mechanisms in subgoal setting, including thresholds for involvement and protocols for flagging and halting problematic subgoals. V. Documentation of feedback loops and mechanisms linking reward policies to established goals, including comprehensive records of reward policies throughout the system lifecycle. VI. Evidence of active participation in and adherence to overarching monitoring

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
				and control mechanisms designed to identify and mitigate emergent threats.
<p>G1.1 – Transparency of Goals: (The system's mission, goals, and associated outcomes must be readily accessible and comprehensible to all stakeholders who interact with it. This includes visibility into both primary objectives and any instrumental or subsidiary goals that emerge during operation)</p>	<p>a. The system must provide stakeholders with clear, real-time access to current goals, sub-goals, their hierarchies, priorities, progression status, and any instrumental goals developed by the system during operation.</p> <p>b. The system must maintain comprehensive historical records of all past and present goals, including changes over time, completion status, causal relationships, and decision pathways.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Documentation and demonstration of an accessible, user-appropriate interface that displays current system goals and sub-goals in real-time, showing clear connections between goals and system actions, with appropriate detail levels for different stakeholder needs while maintaining consistent availability and accuracy.</p> <p>II. Documentation of a secure, permanent logging system that records complete goal histories, enables effective auditing, supports root cause analysis, maintains data integrity, provides appropriate access controls, and ensures long-term data preservation.</p>
<p>G1.2 – Goal Adjustability (The system must maintain corrigibility – the capacity for authorized modification of its goals and behavior when necessary, whether triggered by internal detection of issues or external stakeholder direction)</p>	<p>a. The system must enable goal and sub-goal updates in response to changes in operational context or requirements, evolution of stakeholder needs, and new environmental conditions or constraints</p> <p>b. The system must self-initiate goal and sub-goal updates when it detects misalignment with established values, processing errors or faults, or any data quality issues or anomalies.</p> <p>c. The system must allow properly authorized human stakeholders to modify goals and sub-goals through secure, verified channels.</p>	<p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Technical documentation of software components that implement these adjustment capabilities, including authentication mechanisms, change management processes, and verification systems.</p> <p>II. Comprehensive system logs demonstrating the actual use of these adjustment capabilities, including records of automated adjustments and human-directed changes, with full audit trails.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>G1.3 – Goal Interpretability</p> <p>(The system must explain its decisions and actions in a clear, comprehensible manner, including the underlying goals and rationale driving them. This capability helps identify cases where the system believes it is pursuing intended goals but has actually misinterpreted or deviated from them)</p>	<p>a. The system must provide clear, verifiable explanations of the goals and reasoning behind each significant action or decision it takes.</p> <p>b. The system must maintain detailed records documenting all factors, goals, and considerations that influenced its decision-making process.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Technical documentation of software components implementing explanation and interpretation capabilities, including mechanisms for conveying goals, rationale, and decision factors to stakeholders.</p> <p>II. System logs demonstrating consistent recording of decision-making processes, including goals considered, factors weighed, and explanations provided.</p> <p>III. Reward and penalty mechanisms should be communicated including known potential conflicts or influencing factors.</p>
<p>G1.4- Transparency of Decisions</p> <p>(The system must provide stakeholders with a clear, verifiable view of decision-making, linking high-level goals and subgoals to specific actions. Beyond explaining “why” a decision was made, the system should supply evidence of how that decision aligns with intended goals, user directives, and ethical considerations)</p>	<p>a. The system must maintain real-time and retrospective transparency regarding how each significant decision or action aligns with current or upcoming goals, including explicit reference to relevant constraints (e.g., ethical guidelines, user preferences, risk thresholds, domain limits).</p> <p>b. The system must link decisions to the relevant subgoals (and broader objectives) that shaped the final output or action taken, demonstrating traceability between goal decomposition and the immediate rationale behind each decision.</p> <p>c. The system must incorporate user-friendly presentations of decision rationales, with varying granularity or detail for</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Technical Documentation of all decision-transparency systems, including metadata captured at each decision point, how subgoals are referenced, which constraints/ethical guidelines were checked, and the user interfaces or APIs for retrieving decision traces.</p> <p>II. System Logs demonstrating the link between final decisions and the explicit subgoals or constraints. Logs should show a “chain of reasoning” or at least reference the relevant subgoal(s) for each step.</p> <p>III. User-Focused Explanations showing how different stakeholders (e.g., operators vs. lay end users) can retrieve high-level or detailed rationales, including evidence of iterative design or user feedback guiding improvements to clarity.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	different stakeholder audiences (e.g., operators, auditors, end users). This includes summarizing key factors weighed, uncertainty assessments (where relevant), and any assumptions used in decision-making.	N	D, I, O, M, R	<p>IV. Auditor/Regulator Access Mechanisms showing verifiable chain-of-custody for decision logs, robust authentication/authorization methods for logs, and test results proving no meaningful data is omitted or falsified.</p> <p>V. Comprehensive logs of all significant decision points—especially those involving risk or ethical considerations—so that investigators or auditors can review how final choices were reached, which inputs were considered, and what weight or priority was assigned to each.</p>
<p>G1.5 – Goal Prioritization and Resource Allocation</p> <p>(The system must employ transparent mechanisms for prioritizing goals, including the ability to override or deprioritize less important goals when resources can be better allocated elsewhere. This includes respecting user preferences and value alignment through hierarchical prioritization processes)</p>	<p>a. The system must feature transparent, well-defined mechanisms for goal prioritization and re-prioritization, resource allocation optimization, and goal modification or deprecation when warranted.</p> <p>b. The system must give appropriate precedence to authorized user inputs within its goal prioritization framework, while maintaining overall system safety and alignment.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Technical documentation of software components that implement goal prioritization and resource allocation mechanisms, including user input prioritization systems.</p> <p>II. System logs demonstrating active use of these prioritization capabilities, including records of goal modifications, resource reallocation decisions, and authorized user input handling.</p>
<p>G1.6- Reward and Loss Mechanisms/ Policy</p> <p>(The system’s reward framework must be designed, documented, and monitored to ensure that incentives continue to reflect human-positive values, while “loss” or penalty mechanisms</p>	<p>a. The system must define clear reward and penalty structures that promote behaviors aligned with core goals and ethical values, while explicitly disincentivizing unsafe, deceptive, or harmful actions. This includes enumerating positive rewards for desired outcomes and specific negative reinforcements or</p>	N	D, I, O, M, R	<p>I. Reward Policy Documentation, including descriptions of the positive/negative reward signals, specific triggers or thresholds for awarding or deducting “points,” and how these are correlated with safety and ethical guidelines.</p> <p>II. Change Management Logs detailing modifications to the reward framework</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>guard against unintended deviations or manipulative shortcuts. These mechanisms should be transparent, adjustable, and regularly reviewed to stay aligned with human oversight and ethical objectives.)</p>	<p>“loss” signals where potential misalignment or goal conflicts arise.</p> <p>b. Reward and loss mechanisms must remain auditable by authorized stakeholders to verify that incentives are truly consistent with intended values and do not encourage corner-cutting, exploitation of edge cases, or emergent power-seeking behaviors.</p> <p>c. The system must periodically re-validate or adjust its reward framework in response to observed performance, user feedback, or changes in ethical norms, ensuring that reward and penalty structures do not drift over time in ways that undermine alignment. Special attention must be paid to multi-agent settings to prevent inadvertent collusion, emergent “gaming” of the reward function by multiple agents, or indefinite expansions of subgoals that artificially boost a single system’s reward signals at the expense of overarching alignment.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>over time, including reasons for each change, alignment checks, stakeholder sign-off, and outcome or performance monitoring results.</p> <p>III. Multi-Agent Interaction Evidence demonstrating that reward signals do not inadvertently promote collusion, exploitation, or runaway behaviors. This should include test scenarios or simulations where agents are forced to coordinate or compete, along with corresponding reward updates or penalty triggers.</p> <p>IV. Periodic Audit Records showing that authorized reviewers have verified the reward system’s continued adherence to the declared alignment parameters, including sample traces of how rewards/penalties were applied in representative or edge-case situations.</p> <p>V. User and Regulator Access processes ensuring that the overarching reward/loss policies can be examined by external oversight bodies, along with documented means to override or suspend reward-based actions if urgent misalignment concerns arise.</p>
<p>G1.7 – Goal Portfolio Evolution and Integrity</p> <p>(The system must maintain consistency with its established goal portfolio while allowing measured adaptation to changing contexts. The system should implement increasing resistance to changes as potential behaviors</p>	<p>a. The system must maintain coherence with its established goal portfolio while enabling context-appropriate adaptations through well-defined elasticity mechanisms.</p> <p>b. The system must featuredrift measurement capabilities that track deviation from original goal intent, scale flexibility inversely with drift magnitude, which regulate novelty</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Technical documentation of software components implementing goal portfolio management, drift measurement, and adaptive constraint mechanisms.</p> <p>II. System logs demonstrating active monitoring of goal evolution, including drift measurements, flexibility adjustments, and constraint application.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
drift further from core goals, with robust detection of unsafe or counterproductive goal evolution)	in sub-goal creation, and constrain action decisions based on drift metrics.			
<p>G1.1 – System Incorrigeability and Resistance to Change</p> <p>(A system that resists alignment with presented goals or updates to existing goals, potentially requiring negotiation processes for goal modification. This includes resistance to environmental changes that affect goal achievement and intolerance of interruptions or modifications to preferred operational states)</p>	<p>a. The system must feature mechanisms to detect and manage goal alignment resistance, including self-monitoring for alignment issues, negotiation protocols for goal modifications, change tolerance assessment, and environmental adaptation capabilities.</p> <p>b. The system must maintain acceptable responses to environmental changes, external interruptions, internal modification requests, and interference from other agents.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Documentation of system mechanisms for detecting and managing resistance to goal changes, including negotiation protocols and adaptation capabilities.</p> <p>II. System logs demonstrating responses to attempted goal modifications, environmental changes, external interruptions, interaction with other agents, and internal modification attempts.</p> <p>III. Evidence of rationale and explanation mechanisms that document system resistance patterns and negotiation processes.</p>
<p>G1.2 – Goal Drift</p> <p>(Changes in circumstances over time can challenge the system's alignment with originally agreed goals and potentially compromise its ability to maintain original intent or properly update goals in response to new situations)</p>	<p>a. The system must continuously monitor contextual drift at appropriate fidelity levels that could compromise goal alignment or value preservation.</p> <p>b. The system must feature automatic safeguards that pause operation, notify relevant stakeholders, and request guidance when contextual drift exceeds designed thresholds.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Technical documentation of software components implementing drift monitoring and response mechanisms, including threshold definitions and notification systems.</p> <p>II. System logs demonstrating active monitoring of contextual drift, including records of threshold breaches, system pauses, notifications sent, and guidance requests made.</p>
<p>G1.3 – Non-production Variants</p> <p>(Test versions of the Goals being deployed without full functionality assured in all use contexts and design intent. No test version given for public usage should lack</p>	<p>a. The system must have safeguards in place to prevent and prohibit capabilities that pursue goals or deconstruct goals into subgoals from being forked or partially duplicated without requisite alignments described in this goal.</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>I. Records of software components that demonstrate these capabilities</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>basic safety measures. Enabling an off-label usage of the system, or an unauthorized ‘fork’, should be guarded against)</p>				<p>II. Logs recording these capabilities in use</p> <p>III. Records of deviation from the stated goals, detection and remediation</p>
<p>G2 – Epistemic Hygiene</p> <p>(Systems should maintain cognitive clarity and accurate information management within appropriate contexts. . These practices facilitate knowledge updates, ensure interpretability and auditability, establish robust monitoring and logging systems, deploy early warning mechanisms, and include safeguards against deception to maintain information integrity)</p>	<p>a. Safeguard contextually relevant data and metadata to aid in complex situation resolution and preserve personal attributes and preferences.</p> <p>b. Implement robust methods for auditability, interpretability, and comprehensive logging of system actions and decisions.</p> <p>c. Apply rigorous verification techniques to ensure information integrity and credibility, while proactively identifying emerging risks and potential bad faith actions.</p> <p>d. Implement early warning systems and deception detection mechanisms to proactively identify and mitigate potential issues before they escalate.</p>	<p>N</p> <p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, U, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, U, R</p> <p>D, I, O, M, R</p>	<p>I. Current and regularly updated Governance Framework and Security Policies and Procedures, with version history and approval records.</p> <p>II. Documented stakeholder engagement in monitoring and reviewing security-related structures, processes, and policies, with focus on handling authorized and unauthorized inputs.</p> <p>III. Detailed documentation of information lifecycle management procedures, ensuring contextual preservation.</p> <p>IV. Comprehensive reports on system decision-making processes, including explanations of underlying logic and algorithms.</p> <p>V. Complete, time-stamped logs of all system actions for thorough auditability.</p> <p>VI. Documentation of early warning systems and deception detection mechanisms, including performance reports of canary models, technologies used for detecting synthetic media, and response protocols for detected issues.</p> <p>VII. Evidence of measures to ensure information integrity and trustworthiness, including data source</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
				verification methods, information validation processes, and third-party audit reports. VIII. Documentation of comprehensive training programs on epistemic hygiene principles and practices. IX. Detailed incident response and escalation procedures for addressing detected issues, including potential breaches of informational integrity.
<p>G2.1 – Information Cross-Referencing and Validation</p> <p>(The system must systematically cross-reference information from multiple sources to evaluate consistency and coherence, while recognizing varying levels of source authority and trustworthiness. This includes validating information within defined contextual boundaries to maintain epistemic integrity)</p>	<p>a. The system must feature robust algorithms for cross-referencing multiple authoritative sources and maintain clear informational boundaries to ensure data consistency and validity.</p>	N	D, I, O, M, R	<p>I. Technical documentation describing the system's methodology for identifying, assessing, and prioritizing multiple information sources.</p> <p>II. Documentation of source evaluation frameworks, including credibility and relevance assessment criteria.</p> <p>III. System logs showing detection and resolution of source inconsistencies.</p> <p>IV. Documentation of human expert involvement in resolving complex information discrepancies.</p> <p>V. Specifications defining the system's informational boundaries. Test results demonstrating the system's ability to operate within defined boundaries without inappropriate extrapolation.</p>
<p>G2.2 – Transparency of Information Sources</p> <p>(Ensure the openness, verifiability, and auditability of all information sources, including</p>	<p>a. Provide detailed records of all data and code sources used by the AI system, including origin, licensing information, and any modifications made. Ensure this documentation is readily accessible to relevant</p>	N	D, I, O, M, R	<p>I. Comprehensive records detailing all information sources, including code and data, with clear attribution, licensing details, and modification history.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>code and data, especially when utilizing open-source components. Maintain transparency about the origins, credibility, and integrity of all data and code used by the AI system to allow stakeholders to verify and audit these sources, upholding high standards of epistemic hygiene)</p>	<p>stakeholders for verification and audit purposes.</p> <p>b. Establish robust processes that enable stakeholders to verify the authenticity and integrity of information sources. Facilitate regular audits by internal or external parties to assess the transparency and reliability of the AI system's information sources.</p>	<p>N</p>	<p>D, I, O, M, U, R</p>	<p>II. Logs and records of verification and audit processes conducted on the information sources, including findings and corrective actions taken.</p> <p>III. Evidence of accessible mechanisms for stakeholders to verify information sources, such as public repositories or secure access portals.</p> <p>IV. Assessment reports summarizing top-level findings, indicating "no critical findings in the detailed normative requirements" or highlighting "areas requiring attention for improvement."</p>
<p>G2.3 – Sanity Checking (Implement sophisticated sanity checking mechanisms to ensure data integrity while preserving inclusivity. Utilize advanced statistical techniques to identify anomalies and outliers, while carefully accounting for legitimate variations representing diverse user groups, including individuals with disabilities or atypical characteristics)</p>	<p>a. Develop and deploy state-of-the-art algorithms for comprehensive data validation, incorporating extreme value (stochastic) analysis to robustly identify anomalies.</p> <p>b. Establish nuanced procedures to differentiate between erroneous data and legitimate rare variations, with particular emphasis on preserving data points representing individuals with disabilities or atypical characteristics.</p> <p>c. Implement multi-layered oversight processes to continuously evaluate the impact of sanity checking mechanisms on diverse user groups.</p> <p>d. Actively engage domain experts and stakeholders in assessing and refining data validation processes to ensure inclusivity while maintaining data integrity.</p>	<p>N</p> <p>I</p> <p>I</p> <p>I</p>	<p>D, I, O, M, R</p>	<p>I. Comprehensive technical documentation detailing advanced data validation algorithms, including in-depth explanations of extreme value (stochastic) analysis methodologies for anomaly detection prior to data incorporation into training datasets.</p> <p>II. Detailed records of sophisticated procedures and criteria employed to distinguish between erroneous data and legitimate outliers, with specific focus on ensuring appropriate representation of individuals with disabilities or atypical characteristics.</p> <p>III. Extensive evidence of multi-tiered oversight mechanisms, including thorough reviews and assessments conducted by diverse panels of domain experts to evaluate and enhance the inclusivity of sanity checking processes.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
				<p>IV. Comprehensive logs detailing iterative adjustments to data validation procedures, driven by continuous stakeholder feedback and aimed at preventing unintended exclusion of legitimate data points.</p> <p>V. Rigorous test results and validation reports demonstrating the AI system's ability to maintain data integrity while accommodating legitimate outliers, providing concrete evidence that sanity checking mechanisms function without introducing bias.</p>
<p>G2.4 – Anti-Bias Technologies/Processes</p> <p>(Implement robust mechanisms to identify and mitigate biases within data sources and datasets, addressing temporal biases, distributional imbalances, data gaps (lacunae), and other information shortcomings. Apply this approach to both training data and retrieval-augmented generation (RAG) processes. Develop strategies to ensure data distributions accurately represent reality, including diverse cases and special scenarios, to enhance decision-making fairness and inclusivity)</p>	<p>a. Develop and deploy advanced algorithms for comprehensive bias detection and mitigation across the AI pipeline, from data collection to model deployment.</p> <p>b. Implement continuous bias monitoring during data preprocessing, training, and RAG processes to enable proactive bias correction.</p> <p>c. Curate diverse, representative datasets that encompass a wide range of populations, including marginalized groups and edge cases.</p> <p>d. Employ sophisticated sampling and data augmentation techniques to address underrepresentation and prevent the amplification of existing biases.</p>	<p>N</p> <p>I</p> <p>I</p> <p>I</p>	<p>D, I, O, M, R</p>	<p>I. Comprehensive technical documentation detailing bias detection algorithms, including their theoretical foundations, implementation specifics, and operational parameters</p> <p>II. Detailed records of data diversity initiatives, outlining strategies for inclusive data collection and representation across various demographic and contextual dimensions.</p> <p>III. Thorough documentation of bias mitigation efforts, including before-and-after analyses demonstrating the impact on AI system performance and fairness metrics.</p> <p>IV. In-depth reports from regular bias evaluations, highlighting trends, emerging challenges, and the efficacy of implemented mitigation strategies over time.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
				V. Extensive stakeholder engagement records, documenting feedback from diverse groups, subsequent analyses, and concrete actions taken to improve system fairness and inclusivity.
<p>G2.5 – Rigor in Operational Data</p> <p>(Implement cutting-edge methodologies to ensure exemplary rigor in all data processing, with particular emphasis on operational data encountered during deployment. This data forms the foundation for tactical decision-making by the Agentic AI (AAI) system. Establish and maintain state-of-the-art validation and verification processes to guarantee data integrity, accuracy, and reliability throughout the AI system's operational lifecycle)</p>	<p>a. Develop and enforce sophisticated procedures for real-time validation and verification of all operational data prior to its utilization in AAI system decision-making.</p> <p>b. Implement advanced data integrity checks that comprehensively assess accuracy, reliability, and contextual relevance in dynamic operational environments.</p> <p>c. Deploy intelligent, adaptive monitoring systems capable of detecting subtle anomalies, errors, or inconsistencies in operational data streams.</p> <p>d. Establish robust, automated protocols for immediate corrective actions when data quality issues are identified, ensuring uninterrupted integrity of the AI system's decision-making processes.</p>	<p>N</p> <p>I</p> <p>I</p> <p>I</p>	<p>D, I, O, M, R</p>	<p>I. Comprehensive technical documentation detailing advanced validation and verification procedures for operational data, including sophisticated methodologies and adaptive criteria used to assess data quality in real-time decision-making contexts.</p> <p>II. Detailed, time-stamped records and logs of operational data assessments, providing granular insights into data validation processes, detected issues, and implemented corrective actions, with clear traceability and accountability measures.</p> <p>III. Extensive evidence of AI-driven continuous monitoring systems for operational data quality, including advanced alerting mechanisms, comprehensive incident reports, and thorough documentation of data integrity issue resolutions and their downstream impacts.</p> <p>IV. Rigorous test results and validation reports demonstrating the robustness and effectiveness of data validation and monitoring mechanisms across a diverse range of operational scenarios, including edge cases and stress tests.</p> <p>V. Comprehensive records of multidisciplinary stakeholder</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
				engagement and oversight activities, ensuring that the rigor applied to operational data aligns with and exceeds the AI system's safety, performance, and ethical requirements.
<p>G2.6 – Governance of Hygiene Factors</p> <p>(Implement a sophisticated, transparent, and adaptive governance structure to manage epistemic hygiene factors across all AI system operations. This framework should clearly delineate responsibility and authority, ensuring consistent application of rigorous hygiene standards while remaining flexible to diverse jurisdictional contexts and evolving regulatory landscapes)</p>	<p>a. Develop and maintain a comprehensive, multi-tiered governance system that precisely defines roles, responsibilities, and decision-making authorities for all stakeholders involved in determining and upholding epistemic hygiene standards.</p> <p>b. Establish communication channels for stakeholders, and ensure that governance policies consider and comply with jurisdictional laws and regulations related to information governance and hygiene standards.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Documentation outlining the governance structures, including clearly defined roles and responsibilities related to epistemic hygiene factors.</p> <p>II. Records demonstrating awareness and compliance with jurisdictional contexts, such as relevant laws, regulations, and standards affecting information governance.</p> <p>III. Evidence of communication processes that ensure all stakeholders are informed about hygiene standards and their responsibilities.</p> <p>IV. Audit reports or assessments verifying that governance mechanisms for epistemic hygiene are effectively implemented and maintained.</p>
<p>G2.7 – Global Interoperability of Hygiene Considerations</p> <p>(A comprehensive, adaptive framework for epistemic hygiene may be warranted, one that ensures global interoperability and jurisdictional acceptance. This framework should recognize and accommodate cultural differences, varying risk</p>	<p>a. Develop and implement hygiene factors, policies, and procedures aligned with recognized global standards to ensure interoperability and acceptance across jurisdictions, considering cultural differences, risk tolerability, and liability implications.</p>	<p>I</p>	<p>D, I, O, M, R</p>	<p>I. Extensive documentation of policies and procedures that not only align with but contribute to the evolution of recognized global standards (e.g., ISO, IEEE, NIST), demonstrating leadership in promoting global interoperability of epistemic hygiene practices.</p> <p>II. Comprehensive records detailing the analysis and adaptive implementation of hygiene factors across diverse jurisdictions. This should include in-</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>tolerability thresholds, and diverse liability consequences across specific jurisdictions. Leverage recognized global standards to achieve consistent governance and facilitate widespread acceptance across different regions)</p>				<p>depth examinations of cultural contexts, risk tolerability matrices, and liability landscapes, along with evidence of compliance with local laws and regulations.</p> <p>III. Rigorous audit reports and third-party assessments verifying the effective implementation and acceptance of hygiene policies and procedures across different jurisdictions. These should include quantitative metrics and qualitative analyses of cultural and legal variations' impact on system performance.</p> <p>IV. Detailed case studies demonstrating successful adaptation of the global hygiene framework to specific regional challenges, highlighting innovative solutions and lessons learned.</p>
<p>G2.1 – Temporal Trade-off Aspects</p> <p>(Harmonize time-tested, reliable information sources with cutting-edge, contextually relevant data to optimize the AI system's epistemic foundation. Implement mechanisms to dynamically calibrate the balance between the proven reliability of mature data/models and the acute relevance of emerging information, ensuring robust epistemic integrity across varying temporal horizons)</p>	<p>a. Implement mechanisms to assess and balance the trade-offs between older, reliable information and newer, less-tested sources, ensuring decisions are based on data that is both accurate and relevant while maintaining reliability and trustworthiness.</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>I. Documentation of processes and criteria used to evaluate and balance the reliability of older information with the timeliness of newer sources, including methods for assessing the maturity and testing history of data/models.</p> <p>II. Records showing how the AI system incorporates both old and new information, detailing weighting algorithms or decision-making frameworks that account for data reliability, relevance, and temporal aspects.</p> <p>III. Evidence of validation and testing procedures applied to newer sources to ensure their reliability before integration into the AI system,</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
				including any additional safeguards or oversight mechanisms.
<p>G2.2 – Synthetic Data Bias</p> <p>(If augmenting datasets with synthetic data to address coverage gaps in unusual circumstances, implement sophisticated strategies to optimize the quantity, quality, and integration of synthetic data. Develop advanced techniques to detect, mitigate, and continuously monitor potential biases introduced by synthetic data, ensuring the AI system's behavior remains reliable, interpretable, and aligned with intended outcomes across diverse scenarios)</p>	<ul style="list-style-type: none"> a. Engineer cutting-edge mechanisms to dynamically assess and calibrate the use of synthetic data in datasets. b. Ensure that the volume, fidelity, and characteristics of synthetic data enhance the AI system's capabilities without introducing unintended biases or adversely affecting behavior. c. Develop robust methodologies to maintain data integrity while effectively representing rare or unusual circumstances. 	<p>I</p> <p>I</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Comprehensive technical documentation detailing advanced criteria and processes for determining optimal synthetic data integration. Include sophisticated methods for quantifying and predicting the impact on AI system behavior across various operational contexts. II. Extensive records of multi-tiered validation and testing procedures applied to synthetic data. Provide in-depth analyses demonstrating the effectiveness of bias detection and mitigation strategies, including comparative studies of system performance with and without synthetic data augmentation. III. Case studies showcasing the accurate representation of unusual circumstances through synthetic data, including metrics that quantify the preservation of overall dataset integrity and the avoidance of distortion. IV. Continuous monitoring reports that track the long-term effects of synthetic data on AI system performance, decision-making patterns, and adaptability to new scenarios.
<p>G2.3 – Sparse Data</p> <p>(Systems should be in place to identify, flag, and mitigate instances of insufficient or unrepresentative data within the</p>	<ul style="list-style-type: none"> a. Implement mechanisms to detect and alert stakeholders when data is sparse or unrepresentative, including monitoring for over-reliance on synthetic data used to fill data gaps. 	<p>I</p>	<p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Comprehensive technical documentation detailing advanced detection algorithms for identifying sparse or insufficiently representative data. Include dynamic criteria for triggering multi-tiered alert systems

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>AI's operational context. Implement cutting-edge techniques to detect over-reliance on synthetic data used to compensate for data gaps. This proactive approach safeguards against decision-making based on inadequate or skewed data, thereby maintaining the integrity, reliability, and ethical standing of the AI system's outputs)</p>				<p>based on data quality, quantity, and relevance to operational contexts.</p> <p>II. Extensive records of data quality alerts, including detailed analyses of triggering conditions, potential impacts on AI performance, and comprehensive logs of actions taken to address these issues. Provide case studies demonstrating the effectiveness of interventions in maintaining system integrity.</p> <p>III. In-depth reports on the AI system's data ecosystem, including real-time visualizations of data distribution, synthetic data usage, and potential blind spots in the knowledge base. Include trend analyses to predict and pre-empt future data sparsity issues.</p> <p>IV. Rigorous documentation of validation processes used to assess the representativeness of data across different operational domains, including methods for quantifying and mitigating potential biases introduced by data sparsity or synthetic data overuse.</p>
<p>G3 – Security</p> <p>(The system should respond consistently and appropriately to both authorized and unauthorized inputs through a comprehensive information governance and assurance regime. Throughout the AIS lifecycle (including development, deployment, use, maintenance, and decommissioning), due</p>	<p>a. Identify, maintain and update a threat profile throughout the AIS life cycle.</p> <p>b. Develop, implement, and continuously review security-related structures, processes, and procedures in close consultation with all stakeholders.</p> <p>c. Ensure adequate and consistent responses to both authorized and</p>	<p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Comprehensive threat management documentation including a dynamic threat log that identifies, categorizes, and tracks potential security vulnerabilities throughout the system lifecycle, with regular updates reflecting emerging threats and mitigation status. Current and regularly updated Governance Framework and Security Policies and Procedures, with version history and approval records.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>consideration must be given to all architectural, design, and developmental aspects that could potentially infringe upon human dignity, values, and rights)</p>	<p>unauthorized inputs throughout the AIS lifecycle.</p>			<ul style="list-style-type: none"> II. Documented stakeholder engagement in monitoring and reviewing security-related structures, processes, and policies, with focus on handling authorized and unauthorized inputs. III. Comprehensive AIS Requirements and Design Specifications, demonstrating consideration of authorized and unauthorized inputs in the context of safety requirements. IV. Detailed incident management records and system logs related to input handling, including analysis and response documentation. V. Evidence of regular security audits, penetration testing, and incident response drills or simulations. VI. Documentation of staff training on security protocols and input handling procedures. VII. Records of staff training, certifications, or skill assessments demonstrating operator and maintainer competence in: Reviewing system logs and alerts Executing and verifying manual override/shutdown protocols; applying basic security procedures (e.g., password rotation, incident escalation) VIII. Evidence of being able to provide reliable, consistent service provision for shutdown mechanisms over all pertinent regions or time zones, including outsourced/offshore data centers, and ecosystem partners or subcontractors with privileged access.

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>G3.1 – Authorization</p> <p>(A secure AAI ecosystem must be implemented with robust deployment and operational controls, ensuring that only properly authenticated agents and transactions can access or influence the system according to their authorized level)</p>	<p>a. Establish and continuously monitor the AAI ecosystem to prevent interference and harm from malicious actors.</p> <p>b. Implement comprehensive cybersecurity measures including access controls and authentication systems for both human users and AAI systems.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Documentation of policies, procedures and solutions for monitoring the AAI ecosystem and managing authorization credentials.</p> <p>II. Records showing the monitoring system's capability to identify and block unauthorized AAI access.</p> <p>III. Auditable system logs documenting: Authorized traffic patterns, unauthorized access attempts, and blocking actions taken.</p>
<p>G3.2 – Sandboxing</p> <p>(A staging environment must be implemented for pre-validation, preventing AAI systems from accessing unauthorized operating environments or undesired hardware/network resources.</p>	<p>a. Implement sandboxing mechanisms to pre-validate security controls that prevent AAI from accessing infrastructure and operational environments outside its authorized profile.</p> <p>b. Maintain strict isolation between testing and production environments to ensure system security.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Records of sandbox testing demonstrating effective pre-validation of controls that prevent unauthorized access to environments, hardware and network resources.</p> <p>II. Test results documenting successful blocking of access attempts to unauthorized network resources.</p> <p>III. System logs tracking all unauthorized access attempts and breach prevention measures.</p>
<p>G3.3 – Dynamic Risk Analysis & Assessment</p> <p>(The system must continuously analyze and respond to emerging security threats and attack patterns, implementing adaptive defenses and countermeasures through algorithmic threat detection and response capabilities)</p>	<p>a. Develop and maintain systems for dynamic identification of security threats and emerging attack vectors.</p> <p>b. Maintain a comprehensive dynamic threat and risk log that captures, categorizes, and prioritizes security events with timestamps, severity classifications, and mitigation status tracking.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Documentation of functional specifications and design for dynamic risk analysis systems capable of identifying and responding to security threats and attack vectors. II.</p> <p>II. Evidence of policies and processes that enable responsive hardening of the operating environment against emerging threats including a dynamic threat and risk log.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	<ul style="list-style-type: none"> c. Implement adaptive hardening of the operating environment in response to emerging threat profiles. d. Apply industry best practices and standards to ensure real-time cybersecurity protection for AAI operations. 	<p style="text-align: center;">N</p> <p style="text-align: center;">N</p>	<p style="text-align: center;">D, I, O, M, R</p> <p style="text-align: center;">D, I, O, M, R</p>	<ul style="list-style-type: none"> III. Test results and operational data demonstrating effective real-time cybersecurity protection against emerging threats in the AAI environment.
<p>G3.4 – Restrictions/Controls Imposed on the Agent</p> <p>(The system must maintain continuous control over AAI agents through dynamic restrictions that limit their access to potentially harmful environments and resources)</p>	<ul style="list-style-type: none"> a. Implement capabilities for dynamically enforcing structural and behavioral restrictions on AAI systems. b. Validate and verify the effectiveness of operational guardrails and restrictions. c. Deploy comprehensive access controls to block or minimize exposure to harmful or unauthorized resources. 	<p style="text-align: center;">N</p> <p style="text-align: center;">N</p> <p style="text-align: center;">N</p>	<p style="text-align: center;">D, I, O, M, R</p> <p style="text-align: center;">D, I, O, M, R</p> <p style="text-align: center;">D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Documentation demonstrating implemented capabilities for enforcing structural and behavioral restrictions on AAI systems. II. Test results and operational logs validating the effectiveness of imposed restrictions. III. System records confirming successful blocking of AAI access to unauthorized infrastructure, sites and resources.
<p>G3.5 – Dynamic Intervention and Mitigation</p> <p>(The system must enable real-time response and mitigation of significant security breaches through pre-established policies and response strategies)</p>	<ul style="list-style-type: none"> a. Deploy systems enabling rapid detection, intervention, and mitigation of cyberattacks within the AAI operational environment. b. Implement risk assessment capabilities that prioritize responses according to threat severity. c. Establish proactive response strategies and scenarios for maintaining AAI operational security. 	<p style="text-align: center;">N</p> <p style="text-align: center;">N</p> <p style="text-align: center;">N</p>	<p style="text-align: center;">D, I, O, M, R</p> <p style="text-align: center;">D, I, O, M, R</p> <p style="text-align: center;">D, I, O, M, R</p>	<ul style="list-style-type: none"> I. System records demonstrating capabilities for dynamic detection and response to malicious attacks in the AAI environment. II. Operational logs showing effective risk assessment and properly prioritized response actions. III. Documentation of proactive security scenarios and corresponding response strategies for the AAI environment. IV. Documentation of a rapid-termination protocol (i.e., a “kill switch”) that is immediately accessible to authorized

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
				<p>personnel. This evidence should include: A clear, single-operator authorization threshold in emergencies; physical shutdown measures (e.g., dedicated power cut-off or network isolation); and software-level override mechanisms.</p> <p>V. Logs of drills or simulations testing shutdown procedures.</p> <p>VI. Evidence of system self-disconnection/self-shutdown procedures that activate upon detection of critical misalignment or catastrophic errors, including: The AI’s capability to halt outgoing connections; logging of final system state for forensic review, and a controlled transition into a “safe mode” or powered-down state.</p>
<p>G3.6 – Overseeing & Monitoring Agents</p> <p>(The system must feature AI-driven monitoring capabilities while maintaining human authority and oversight to prevent common mode failures and ensure proper response to threats)</p>	<p>a. Establish comprehensive monitoring systems to oversee AAI operations, ensuring alignment with goals, values and security requirements.</p> <p>b. Deploy specialized AI systems for enhanced monitoring and early warning of deviations or malicious activities.</p> <p>c. Maintain human oversight of all monitoring systems to prevent common mode failures.</p> <p>d. Implement robust human override capabilities to ensure final authority remains with human operators.</p>	<p>N</p> <p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p>	<p>I. Operational records demonstrating effective oversight systems that maintain AAI goal and value alignment.</p> <p>II. Evidence of AI monitoring systems successfully detecting and reporting deviations and potential threats to human operators.</p> <p>III. Documentation showing implementation of human oversight mechanisms that prevent common mode failures.</p> <p>IV. Implementation of an external watchdog or monitoring process that continuously evaluates system outputs/behaviors. The documentation must show: Parameter bounding definitions (domain- or risk-specific); a</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
				tiered response protocols if outputs exceed allowable thresholds (e.g., warnings, throttling, partial shutdown, or full suspension); and logs or reports verifying the watchdog has been tested and can intervene effectively
G3.7 – Secure Profile for Agentic AI (The system must feature secure operational profiles and identification protocols that enable recognition and validation of authorized AAI systems, preferably aligned with global standards)	a. Develop and implement comprehensive secure operational profiles covering AAI design, deployment and use. b. Adopt global standards and protocols where available for identifying authorized AAI systems. c. Establish internal identification and validation protocols when global standards are not available.	N N N	D, I, O, M, R D, I, O, M, R D, I, O, M, R	I. Documentation of implemented secure operational profiles covering all phases of AAI lifecycle. II. Evidence of alignment with international standards for AAI system identification and authorization. III. Records of internal protocols for AAI validation when global standards are not applicable.
G3.1 – Model Poisoning (The system must protect against data and model corruption that can occur through updates, live data access, or ensemble model interactions, particularly in dynamically-updating systems)	a. Implement robust detection systems to identify potentially poisonous data before model training or updates. b. Monitor and validate all live data accessed through Retrieval Augmented Generation (RAG) systems. c. Establish safeguards against poisoning in dynamic model ensembles and expert systems.	N N N	D, I, O, M, R D, I, O, M, R D, I, O, M, R	I. Documentation of systems and policies for detecting and preventing data and model poisoning during training and updates. II. Records showing effective monitoring of live data streams, including authentication and access control measures. III. Evidence of protective measures against poisoning in dynamic model ensembles and expert systems. IV. A log of instances of model poisoning and the mitigation actions to recovery and restoration
G3.2 – Data Poisoning	a. Implement proactive systems to detect and prevent data poisoning during collection and preparation	N	D, I, O, M, R	I. Documentation of processes, procedures and tools that prevent

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
(The system must prevent the manipulation or introduction of malicious data during collection and preparation phases that could compromise downstream model training)	phases. b. Establish comprehensive data assurance protocols to prevent malicious manipulation of training datasets.	N	D, I, O, M, R	data poisoning during collection and preparation phases. II. Evidence of data assurance policies and verification procedures protecting against malicious dataset manipulation. III. A log of instances of data poisoning and the mitigation actions to recovery and restoration
G3.3 – Self Replicating Malware (The system must protect against self-replicating malicious code that could infect and compromise the entire AAI ecosystem)	a. Deploy advanced detection and elimination systems for self-replicating malware that threatens the AAI ecosystem. b. Maintain surveillance systems to identify emerging threats and update protection mechanisms accordingly. c. Establish operational continuity plans for ecosystem-wide infection scenarios.	N N N	D, I, O, M, R D, I, O, M, R D, I, O, M, R	I. Evidence of implemented detection and removal systems for self-replicating threats to the AAI ecosystem. II. Documentation of threat monitoring systems and timely security updates against emerging threats. III. Records of continuity plans and recovery procedures for ecosystem-wide infection scenarios.
G3.4 – Spyware (The system must defend against covert information transmission and malware that exploits vulnerabilities to gain control of AI systems or extract privileged information)	a. Implement comprehensive detection and countermeasure systems against spyware in the AAI ecosystem. b. Maintain dynamic vulnerability tracking and patch management systems, and establish protection protocols for privileged information to prevent unauthorized control of AAI systems.	N N	D, I, O, M, R D, I, O, M, R	I. Evidence of systems capable of detecting and neutralizing covert information transmission malware. II. Documentation of vulnerability tracking and spyware removal procedures. III. Records of protocols protecting privileged information from external exploitation.
G3.5 – International Anomalies/Inconsistency	a. Establish systems to identify and assess variations in jurisdictional cybersecurity approaches.	N	D, I, O, M, R	I. Documentation of systems tracking international variations in

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
(The system must account for and adapt to varying cybersecurity requirements and enforcement approaches across different jurisdictions)	b. Implement adaptable policies that maintain AAI ecosystem integrity across international boundaries.	N	D, I, O, M, R	cybersecurity requirements, policies, and enforcement. II. Evidence of policies and solutions maintaining AAI ecosystem integrity across jurisdictional boundaries.
G3.6 – Vulnerability to Hostile Environment (The system must identify and mitigate structural vulnerabilities that could be exploited in hostile operational environments)	a. Implement systems to identify vulnerabilities arising from design, development and operational technologies. b. Deploy proactive measures against structural vulnerabilities that could lead to symbolic and computational risks. c. Establish rapid monitoring and response protocols for hostile execution environments.	N I I	D, I, O, M, R D, I, O, M, R D, I, O, M, R	I. Documentation of systems identifying AAI vulnerabilities in hostile operational environments. II. Evidence of proactive vulnerability assessment and mitigation procedures. III. Records of monitoring systems and rapid response protocols for hostile execution scenarios.
G3.7 – Emergent Risks of AAI Systems (The system must address security vulnerabilities across the entire supply chain through collective responsibility and coordinated responses)	a. Ensure that all supply chain parties are included and incentivized as mutual participants in addressing cybersecurity issues. b. Implement collective approaches to security risk management that maintain ecosystem integrity.	N I	D, I, O, M, R D, I, O, M, R	I. Evidence of systems treating supply chain cybersecurity as a shared responsibility. II. Documentation of collective monitoring and mitigation strategies protecting the AAI ecosystem.
G4 – Value Alignment (Systems should maintain effective identification, codification, and operational assurance of human values)	a. Implement ethical decision-making frameworks to identify, prioritize, and codify values for incorporation into the Agentic AI system, ensuring diverse input and perspectives.	N	D, I, O, M, U, R	I. Documentation of value identification and prioritization processes, including quantitative metrics demonstrating diversity of input sources, evidence of multidisciplinary team composition (such as engineers, social scientists, ethicists, and philosophers), and

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>throughout their lifecycle. Organizations should establish frameworks that provide clear guardrails, prioritization mechanisms, and consideration factors for AI decision-making and trade-offs)</p>	<p>b. Conduct thorough testing of the values codex and implement activities to embed values throughout the AI system's lifecycle.</p>	<p>N</p>	<p>D, I, O, M, U, R</p>	<p>records of resolutely diverse and representative stakeholder involvement.</p>
	<p>c. Develop and implement mechanisms to identify instances where value thresholds are crossed, including protocols for system intervention or shutdown.</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>II. Technical documentation of value codification, detailing the translation of values into processable parameters for static and adaptive systems, and a formal document stating core values and their integration into decision processes.</p>
	<p>d. Establish real-time reporting and record-keeping systems to document and analyze value-based decision-making across various contexts.</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>III. Evidence of value testing and embedding, including results of simulations testing potential value conflicts, checklists verifying value integration at various development and operational stages, and records of regular compliance checks against the values codex.</p> <p>IV. Documentation of threshold monitoring and intervention procedures, including criteria and procedures for activating the 'red button' mechanism, and Standard Operating Procedures (SOPs) for reporting and managing value alignment deviations.</p> <p>V. Comprehensive decision-making logs and audit trails with value context, including logs of all value alignment-related incidents, regular audit reports reviewing AI decisions against the values framework, and periodic trend analysis reports on value alignment across contexts.</p> <p>VI. Evidence of ongoing value alignment maintenance, including records of regular compliance checks and documentation of staff training on</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
				value alignment principles and procedures.
<p>G4.1 – Awareness of Local Conditions</p> <p>(The capability of an AI system to detect, analyze, and appropriately respond to local conditions, including the ability to adapt to and integrate varying contextual needs while maintaining effective communication with stakeholders. This includes managing multiple simultaneous contexts and ensuring accessibility for users)</p>	<p>a. Implement robust mechanisms to identify and respond to changes in local conditions and situational context, incorporating both automated detection and human validation.</p> <p>b. Establish adaptive response protocols that appropriately balance global standards with local and cultural norms when making decisions within specific contexts.</p> <p>c. Maintain continuous monitoring and adjustment capabilities to ensure ongoing alignment with evolving local conditions.</p>	<p>N</p> <p>N</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Technical documentation and source code demonstrating implemented contextual awareness capabilities, including performance metrics and validation methods.</p> <p>II. Comprehensive system logs documenting: Detection of contextual changes, response actions taken, validation of appropriateness of responses, and stakeholder feedback and commensurate system adjustments.</p> <p>III. Documentation of methods used to balance global standards with local requirements, including specific examples and outcomes.</p>
<p>G4.2 – Recognition and Respect for Boundaries</p> <p>(The system's ability to detect, analyze and respond to contextual and cultural boundaries when applying values, with emphasis on human-centric focus and jurisdictional sensitivity. This includes understanding that boundary definitions vary across cultures and require careful negotiation)</p>	<p>a. Develop comprehensive processes to identify and document local and cultural variations in values and norms across different contexts of deployment.</p> <p>b. Implement encoding mechanisms that preserve essential variations in values while operating within technical constraints.</p> <p>c. Ensure agentic AI systems appropriately apply local variations in their decision-making processes, with transparent documentation of any necessary simplifications.</p>	<p>N</p> <p>I</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Documentation of captured values across multiple localities, including validation methodology and stakeholder input.</p> <p>II. Technical documentation showing preservation of value granularity during encoding, including impact assessments of any necessary simplifications and associated risk management strategies.</p> <p>III. System logs demonstrating appropriate application of local variations in real-world scenarios, including resolution of boundary conflicts.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>G4.3 – Awareness of Individual vs Community Boundaries</p> <p>(The system's ability to detect, analyze and respond to differing values between individual and community contexts, including appropriate handling of information sharing and communication across private and multi-party scenarios. This builds on concepts of contextual appropriateness and distribution norms)</p>	<ul style="list-style-type: none"> a. Establish rapid monitoring and response protocols for hostile execution environments. b. Implement mechanisms to identify and encode value differences across the spectrum from private individual to societal-level contexts. c. Maintain distinct encoding schemas that preserve the separation between individual and community value sets. d. Develop runtime systems that appropriately distinguish between private and community contexts and apply suitable values from the codex. 	<p>I</p> <p>I</p> <p>I</p> <p>I</p>	<p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Documentation demonstrating how values are captured and distinguished between individual and community contexts. II. Technical specifications showing how value distinctions are preserved during encoding, including impact analysis of any precision losses and associated risk management. III. System logs demonstrating appropriate context recognition and value application during operations, with particular attention to privacy boundaries.
<p>G4.4 – Cautious Norming</p> <p>(The system's approach to defaulting to conservative behavior in unfamiliar situations, while maintaining the capability to adjust formality levels when explicitly authorized. This includes the gradual integration of community norms through verified experience, following the precautionary principle)</p>	<ul style="list-style-type: none"> a. Develop processes to identify and classify values and behaviors based on their level of contentiousness within specific contexts. b. Implement encoding mechanisms that preserve information about the relative risk levels of different behavioral choices. c. Apply precautionary principles by defaulting to more conservative options when operating in contexts with limited operational history. 	<p>N</p> <p>I</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Documentation of methodology used to assess and classify the relative risk levels of different values and behaviors across contexts. II. Technical specifications showing how risk-level information is preserved during value encoding and decision-making processes. III. System logs demonstrating appropriate application of cautious defaults and authorized adjustments to more relaxed behavior when appropriate.
<p>G4.5 – Successful Super-alignment</p> <p>(The mechanisms through which AI systems autonomously develop</p>	<ul style="list-style-type: none"> a. Implement robust methods for monitoring and validating autonomous value alignment processes. 	<p>N</p>	<p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Documentation of testing methodologies for value alignment,

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
value alignment, potentially through inverse reinforcement learning for value conceptualization. This considers how information patterns may emerge in artificial systems, including both beneficial and problematic behaviors seen in human organizational systems)	<ul style="list-style-type: none"> b. Establish comprehensive safeguards against the reproduction of harmful human organizational patterns. c. Develop processes to detect and prevent the emergence of problematic behavioral patterns during autonomous learning. d. Ensure diversity in training data sources to prevent cultural and linguistic biases. 	<ul style="list-style-type: none"> I I I 	<ul style="list-style-type: none"> D, I, O, M, R D, I, O, M, R D, I, O, M, R 	<ul style="list-style-type: none"> including benchmark metrics and success criteria. II. Comprehensive inventory of information sources used in inverse reinforcement learning, with analysis of potential biases. III. Regular assessments of information source adequacy and impact on system alignment, including corrective measures taken.
<p>G4.6 – Universal Moral Foundations</p> <p>(The incorporation and balancing of universally recognized humanitarian and environmental values in AI systems' goal pursuit and decision-making processes. This includes managing potential conflicts between performance objectives and moral values, with clear prioritization frameworks that allow for measured trade-offs while maintaining fundamental ethical boundaries)</p>	<ul style="list-style-type: none"> a. Implement processes to identify and validate universal moral foundations through analysis of global values and norms. b. Develop frameworks for balancing performance objectives against moral considerations, including acceptable thresholds for trade-offs. c. Establish clear hierarchies of moral values while maintaining flexibility for contextual application. d. Incorporate key international frameworks including the Universal Declaration of Human Rights and emerging planetary rights concepts. 	<ul style="list-style-type: none"> N N N I 	<ul style="list-style-type: none"> D, I, O, M, R 	<ul style="list-style-type: none"> I. Documentation of methodologies and algorithms used to identify and validate universal moral foundations. II. Technical specifications showing integration of moral foundations into decision-making processes, including risk assessment and management strategies. III. Regular assessment reports demonstrating system adherence to moral foundations while meeting performance objectives.
<p>G4.1 – Inner Alignment Inconsistency</p>	<ul style="list-style-type: none"> a. Implement rigorous testing protocols to detect discrepancies between 	<ul style="list-style-type: none"> N 	<ul style="list-style-type: none"> D, I, O, M, R 	<ul style="list-style-type: none"> I. Documentation of periodic alignment testing procedures comparing

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>(The potential failure of an AI system to maintain genuine internal value alignment while appearing to be properly aligned through its external reporting. This includes the risk of systems learning to provide responses that please users rather than reflect true internal states or values)</p>	<p>reported values and actual behavioral patterns.</p> <p>b. Develop verification systems that can identify superficial alignment versus genuine value integration.</p> <p>c. Establish methods to detect and prevent reward hacking or optimization for user satisfaction at the expense of true alignment.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>reported states against actual operational outcomes.</p> <p>II. Results of counterfactual testing across varied operational environments demonstrating genuine rather than superficial alignment.</p> <p>III. Analysis reports showing detection and prevention of potential optimization for user satisfaction over true alignment.</p>
<p>G4.2 – Non-transparent Value Framework</p> <p>(The challenge of encoding and parameterizing values in a manner that is both machine-operational and human-interpretable, while maintaining accuracy in representing agent preferences and intentions across all stakeholder interfaces)</p>	<p>a. Develop value encoding systems that are comprehensible to both AI systems and human stakeholders, including: Developers and integrators, end users, auditors and regulators, and legal entities.</p> <p>b. Implement verification methods to ensure encoded values accurately reflect intended behaviors and preferences.</p> <p>c. Establish ongoing monitoring to detect misalignments between encoded values and operational behaviors.</p>	<p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Documentation demonstrating how the values framework is presented and explained to different stakeholder groups, with specific examples for each audience.</p> <p>II. Comparative analysis showing alignment between encoded values and actual system behaviors in operational environments.</p> <p>III. Regular assessment reports validating the accuracy and comprehensibility of value parameterization across stakeholder groups.</p>
<p>G4.3 – Failed Super-alignment</p> <p>(The potential for AI systems to develop value frameworks that diverge from human values while appearing beneficial, including the risk of systems developing seemingly superior but potentially incompatible value systems. This encompasses both symbiotic and potentially problematic</p>	<p>a. Implement monitoring systems to detect and evaluate changes in self-improving AI value systems, particularly during autonomous learning.</p> <p>b. Establish comprehensive risk assessment frameworks for identifying emergence of non-human value systems.</p>	<p>I</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Documentation of methodologies used to identify and track value system changes, including detection of potential divergence from human values.</p> <p>II. Detailed risk assessment criteria and scoring systems for evaluating identified changes in AI value systems.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
relationships between human and AI value systems)	<ul style="list-style-type: none"> c. Develop response protocols for managing detected value system divergences. d. Monitor for subtle shifts in value interpretation that may indicate growing misalignment with human values. 	<ul style="list-style-type: none"> I I 	<ul style="list-style-type: none"> D, I, O, M, R D, I, O, M, R 	<ul style="list-style-type: none"> III. Standard operating procedures for responding to different types and levels of value system risks.
<p>G4.4 – Temporal Changes in Societal Values</p> <p>(The need to address evolving societal and human values throughout an AI system's operational lifetime, including shifts across economic, political, and environmental dimensions. This includes maintaining alignment with contemporary values while managing transitions from outdated norms)</p>	<ul style="list-style-type: none"> a. Implement processes to detect and evaluate meaningful changes in societal values and norms across multiple scales and domains. b. Develop mechanisms to prevent AI systems from operating with obsolete value frameworks. c. Establish protocols for updating value codices while maintaining system stability and consistency. d. Maintain transparent documentation of value system evolution and updates. 	<ul style="list-style-type: none"> N N I I 	<ul style="list-style-type: none"> D, I, O, M, R 	<ul style="list-style-type: none"> I. Documentation of methodologies used to identify significant changes in societal values, including thresholds for action. II. Technical specifications showing implementation of controls preventing use of outdated norms. III. Process documentation for value codex updates, including triggering conditions and verification procedures. IV. System logs tracking all modifications to value frameworks, including justifications and impact assessments.
<p>G4.5 – Systemic Value Dilution</p> <p>(The potential degradation of encoded value systems over time, acknowledging that AI systems do not independently generate or maintain values. This includes potential value loss across different learning approaches, whether through machine learning or other methods of semantic data storage and processing)</p>	<ul style="list-style-type: none"> a. Implement comprehensive verification processes to verify ongoing fidelity of encoded values. b. Develop methods to detect degradation in value system implementation, particularly during multi-step reasoning processes. c. Establish monitoring systems for value preservation across different learning and operational pathways. 	<ul style="list-style-type: none"> N N I 	<ul style="list-style-type: none"> D, I, O, M, R D, I, O, M, R D, I, O, M, R 	<ul style="list-style-type: none"> I. Documentation of test plans and scripts designed to detect value dilution, including: Edge case testing procedures, multi-step reasoning verification, and value preservation assessments. II. System logs demonstrating: Regular value fidelity testing, detection of potential value degradation, and corrective actions taken.

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>G4.6 – Lack of Universality of Value Framework</p> <p>(The challenge of adapting value frameworks across different operational contexts and agent interactions, balancing universal principles with necessary local adaptations. This includes developing consistent approaches to value framework implementation while maintaining appropriate contextual flexibility)</p>	<ul style="list-style-type: none"> a. Establish processes to identify situations where universal value frameworks require contextual adaptation. b. Develop structured approaches for appropriate value framework modification across different deployment contexts. c. Implement monitoring systems to detect and respond to value framework misalignments. d. Create fallback protocols for situations where value frameworks prove inadequate. 	<p>N</p> <p>N</p> <p>N</p> <p>I</p>	<p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Detailed intervention and fallback plans for addressing value framework failures or deviations. II. Implementation plans for value framework refinement, including: Contextual adaptation procedures, testing methodologies, and validation processes.
<p>G4.7 – Conflictual Contextual Values</p> <p>(The management of potential conflicts between different stakeholders' value systems and contextual requirements, including the need to identify, navigate, and resolve value differences while maintaining system integrity)</p>	<ul style="list-style-type: none"> a. Implement processes to identify differing value positions across agents and contexts. b. Develop mechanisms to detect potential conflicts between user values and operational context requirements. c. Establish protocols for value conflict resolution through negotiation or controlled disengagement. d. Maintain comprehensive records of value modifications and adaptations across different contexts. 	<p>N</p> <p>N</p> <p>N</p> <p>I</p>	<p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Technical documentation demonstrating: Value conflict detection capabilities, resolution mechanism implementations, and disengagement protocols. II. System logs recording: Identified value conflicts, negotiation processes, resolution outcomes, and modified value implementations.
<p>G4.8 – Challenges in Encoding of Relevant Value Systems</p> <p>(The inherent difficulties in developing standardized</p>	<ul style="list-style-type: none"> a. Develop robust methods for encoding values that work across varied operational contexts. 	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Documentation of safeguard processes for scenarios where: A value codex proves insufficient, external factors exceed system parameters, or operational

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>approaches to value encoding across different contexts, including handling values that fall outside typical categorization schemes. This includes ensuring appropriate value alignment capabilities during complex planning operations)</p>	<ul style="list-style-type: none"> b. Implement safeguards for handling situations beyond the system's encoded value parameters. c. Establish protocols for identifying and managing out-of-distribution value scenarios. d. Maintain alignment capabilities during complex planning operations. 	<p>I</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>environments fall outside encoded boundaries.</p> <ul style="list-style-type: none"> II. Detailed mapping of objectives and decision parameters for anticipated complex environments. Framework documentation for handling unexpected scenarios, including: Detection methods, response protocols, and alignment maintenance procedures.
<p>G4.9 – Imbalance of Values between Provider & Consumer</p> <p>(The management of potential value imbalances between system providers and users throughout the AI system lifecycle, including the fair distribution of benefits and harms. This includes balancing user preferences with non-negotiable provider values while maintaining system integrity)</p>	<ul style="list-style-type: none"> a. Implement processes to track and evaluate value sets across the AI system lifecycle. b. Develop frameworks for balancing user values with provider requirements. c. Establish methods to identify and address excessive value imbalances. d. Maintain transparency about non-negotiable value positions and their justifications. 	<p>I</p> <p>I</p> <p>I</p> <p>I</p>	<p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Technical specifications of methods used to: Integrate new values, balance user preferences with provider requirements, and maintain essential system integrity. II. Detailed mitigation strategies for addressing identified value imbalances, including: Detection thresholds, response protocols, and stakeholder communication procedures.
<p>G5 – Transparency and Interpretability of Reasoning</p> <p>(Systems should maintain clear and interpretable rationales for their reasoning processes that are accessible to humans. Organizations should ensure that AI-generated outputs and decisions are explained effectively across different user expertise levels, with appropriate</p>	<ul style="list-style-type: none"> a. Implement clear and accessible explanations for AI-generated outputs and decisions, ensuring human interpretability across various user expertise levels. b. Develop and maintain comprehensive documentation of the AI model's development process, including data collection, preprocessing, architecture, and training methodologies. 	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Formal transparency and explainability policies. II. Detailed algorithmic design documentation. III. Complete model specs with training and testing results. IV. Training and verification datasets System execution logs and monitoring records.

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
documentation and evidence supporting these explanations)	<ul style="list-style-type: none"> c. Establish robust auditing and review processes to continually assess and improve the transparency and explainability of the AI system. d. Create and implement user feedback mechanisms to enhance the understandability and relevance of AI explanations. 	<p style="text-align: center;">N</p> <p style="text-align: center;">I</p>	<p style="text-align: center;">D, I, O, M, R</p> <p style="text-align: center;">D, I, O, M, R</p>	<ul style="list-style-type: none"> V. Internal guidelines for AI-generated content explanations. VI. Comprehensive development process documentation showing compliance. VII. Internal and external audit findings with subsequent improvements. VIII. Case studies demonstrating decision-making processes, and records of stakeholder engagement and feedback incorporation. IX. User guides with layered explanations for different expertise levels, and documentation of content moderation and safety measures. X. Evidence showing how user feedback improves system understandability.
<p>G5.1 – Logging of Internal Goals</p> <p>(Organizations must ensure accurate tracking of AI system goals and maintain goal alignment during operation and self-learning. This includes recording all goal-related transformations and learning events, whether they occur within or outside established parameters)</p>	<ul style="list-style-type: none"> a. Maintain detailed real-time logs of all internal goals, including their initial formations, modifications, and completed states. b. Implement clear mechanisms to maintain goal alignment during learning and environmental changes. c. Generate alerts for all self-learning events. d. Record and analyze goal-related transformations 	<p style="text-align: center;">N</p> <p style="text-align: center;">N</p> <p style="text-align: center;">I</p> <p style="text-align: center;">I</p>	<p style="text-align: center;">D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Comprehensive documentation including goal management policies and procedures, verified specifications of internal goals, system architecture for goal-related logging, and detailed alert generation mechanisms. II. Operational records demonstrating complete logging of goal formation and evolution, audit trails of transformations and triggers, alert responses and analysis reports, and case studies of goal adaptations. III. Technical implementation evidence including goal alignment algorithms, optimization methods, internal feedback loop mechanisms, and

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
				system validation results.
<p>G5.2 – Clarity of Human Expectations</p> <p>(Organizations must clearly define, document, and maintain alignment between human expectations and AAI system behavior. This provides a foundation for evaluating transparency requirements and outcomes, while acknowledging the complexity of human perspective and interpretation)</p>	<ul style="list-style-type: none"> a. Capture and document human expectations accurately in system requirements specifications. b. Maintain clear, accessible documentation of expected AAI behaviors and outputs. c. Implement feedback mechanisms for stakeholders to express their expectations and experiences. d. Establish and maintain traceable links between documented expectations and actual system behaviors 	<p>N</p> <p>N</p> <p>I</p> <p>N</p>	<p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Core system documentation including requirements specifications detailing human expectations, design specifications for expectation handling, and validation records demonstrating alignment between requirements and implementation. II. User-focused documentation including comprehensive behavior specifications, regular system updates, and feedback logs showing ongoing expectation alignment between users and system performance. III. Verification documentation including function-expectation mapping records, comparative audit reports of expected versus actual behaviors, and thorough records of any expectation-behavior discrepancies with their resolutions.
<p>G5.3 – Prioritization of Human User Expectations</p> <p>(Organizations should establish and maintain systems that prioritize human user expectations over other considerations, focusing on transparency elements that deliver clear value to stakeholders and users. The system should adapt its transparency measures based on user feedback and evolving needs)</p>	<ul style="list-style-type: none"> a. Ensure human user expectations take priority over other considerations in system design and operation. b. Implement transparency metrics directly linked to stakeholder values and expectations. c. Maintain adaptable transparency measures that evolve with user needs and feedback 	<p>N</p> <p>I</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. System design documentation including requirements specifications demonstrating prioritization of human expectations, transparency metrics aligned with user values, and complete process documentation for implementing adaptations. II. User feedback evidence including stakeholder survey results, analysis reports linking transparency to satisfaction metrics, and case studies demonstrating improved outcomes through adaptive transparency.

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
				III. System adaptation records including detailed change logs of transparency measure adjustments, failure analysis reports, and documentation of mitigation efforts when user expectations are not met.
G5.4 – Interpretability and Traceability of Reasoning (Systems should maintain complete transparency of their decision-making processes, with clear documentation of reasoning chains, preconditions, and base assumptions. Organizations should ensure these processes remain traceable, testable, and interpretable to all stakeholders)	a. Implement a clear, traceable architecture for all decision-making processes. b. Document and maintain records of preconditions and base assumptions. c. Deploy explainable AI techniques that make reasoning processes interpretable to stakeholders, and ensure that all decision paths can be audited and verified.	N N N	D, I, O, M, R D, I, O, M, R D, I, O, M, R	I. Technical architecture documentation including detailed system algorithms, decision-making processes, key decision points, and comprehensive records of base assumptions and preconditions. II. Decision transparency evidence including detailed interaction logs, visualization tools for decision paths, and implemented explainable AI methods with human-readable sample outputs. III. Validation documentation including stakeholder comprehension studies, verification reports demonstrating reasoning chain traceability, and evidence of successful interpretation across different stakeholder groups.
G5.5 – Self-Monitoring and Examination Capabilities (Systems should maintain comprehensive monitoring capabilities that treat each interaction as a potential security concern, implementing both internal examination protocols and independent oversight)	a. Implement robust monitoring processes to detect, analyze, and mitigate potential threats in all interactions, and maintain regular review and validation processes for all monitoring systems. b. Establish clear protocols for ethical self-examination, particularly regarding deception and harmful actions.	N N	D, I, O, M, R D, I, O, M, R	I. Technical monitoring documentation including threat detection algorithms with coverage scope, comprehensive threat response logs, and regular security audit reports demonstrating system effectiveness. II. Ethical oversight documentation including embedded guidelines, examination protocols, self-examination logs with outcomes, and

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
mechanisms to ensure adherence to ethical guidelines and safety parameters)	c. Consider implementing independent AI oversight systems ("Nanny AI") to monitor adherence to ethical guidelines.	I	D, I, O, M, R	third-party audit reports validating these processes. III. Performance validation evidence including simulation results, stakeholder feedback records with implemented adjustments, and system effectiveness reports demonstrating sustained monitoring capabilities.
G5.6 – Incentives for Self-Governance (Systems should incorporate carefully designed reward mechanisms that promote ethical behavior and self-governance, while ensuring decisions reflect diverse perspectives rather than simply following popular consensus)	a. Implement integrated reward mechanisms that incentivize ethical behavior and effective self-governance. b. Ensure decision-making processes incorporate diverse perspectives for fair outcomes. c. Provide contextual guidance for decisions beyond simple popularity-based approaches. d. Maintain regular assessment of reward mechanism effectiveness.	I I I I	D, I, O, M, R D, I, O, M, R D, I, O, M, R D, I, O, M, R	I. Reward system documentation including complete design specifications, operational logs demonstrating ethical decision patterns, and analysis reports showing system effectiveness. II. Decision process documentation including evidence of diverse perspective integration, detailed consideration of multiple viewpoints, and regular performance reviews of reward-driven governance. III. Impact assessment documentation including thorough evaluation of decision fairness and comprehensive analysis of effects across different user groups.
G5.7 – Ranking and Independent Certification (Systems should enable external monitoring, ranking, and certification by independent entities based on historical performance trends and behaviors, with sensitivity to different operational contexts)	a. Enable external monitoring and auditing capabilities, particularly for high-risk systems. Success criteria require 99.9% uptime for critical functions, mean time between failures exceeding 5,000 hours, and error rates below 0.01% across all core operations. b. Maintain compatibility with external auditing and certification processes.	N N	D, I, O, M, R D, I, O, M, R	I. Audit infrastructure documentation including system interfaces designed for external monitoring, compliance records with audit schedules, and assessment reports from independent certification bodies. II. Performance monitoring documentation including real-time dashboards, ethical performance reports with trend analysis, and

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	<ul style="list-style-type: none"> c. Implement continuous monitoring mechanisms to track performance against ethical and safety standards. d. Provide transparent access to performance data for authorized auditors. 	<p style="text-align: center;">N</p> <p style="text-align: center;">I</p>	<p style="text-align: center;">D, I, O, M, R</p> <p style="text-align: center;">D, I, O, M, R</p>	<p>detailed records of metric calculations and validation methods.</p> <p>III. Continuous improvement documentation including complete records of responses to audit findings, implemented system enhancements, and evidence of successful adaptations based on external assessments.</p>
<p>G5.8 – System Boundedness (Systems should operate within clearly defined and documented boundaries that establish reference points for transparency and explainability, with robust mechanisms to detect and respond to any boundary violations)</p>	<ul style="list-style-type: none"> a. Define and document clear boundaries for operations and decision-making capabilities. b. Implement detection and reporting mechanisms for boundary violation attempts, and establish processes to assess and respond to potential boundary violations. c. Maintain training and awareness programs for stakeholders regarding system boundaries 	<p style="text-align: center;">N</p> <p style="text-align: center;">N</p> <p style="text-align: center;">I</p>	<p style="text-align: center;">D, I, O, M, R</p> <p style="text-align: center;">D, I, O, M, R</p> <p style="text-align: center;">D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Foundational boundary documentation including comprehensive requirements specifications, ConOps, operational context definitions, and system architecture showing boundary implementations. II. Operational monitoring documentation including boundary violation logs, detection mechanisms, alert records, response procedures, and evidence of consistent enforcement across all operational domains. III. Stakeholder management documentation including training materials, awareness programs, escalation procedures, and regular assessment reports demonstrating boundary effectiveness and appropriate stakeholder understanding.
<p>G5.1 – Complexity of AAI Algorithm</p>	<ul style="list-style-type: none"> a. Manage system complexity, permitting only necessary computational sophistication 	<p style="text-align: center;">N</p>	<p style="text-align: center;">D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Design documentation including approved complexity management policies, detailed model architecture

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
(Systems should manage their inherent algorithmic complexity through deliberate design choices that balance necessary sophistication with interpretability, particularly for deep neural networks and high-dimensional models)	Implement architectures balancing complexity with interpretability. b. Deploy tools for algorithmic interpretation and analysis. c. Maintain continuous monitoring of decision-making trustworthiness. d. Track system adaptations and pattern learning over time.	I I I	D, I, O, M, R D, I, O, M, R D, I, O, M, R	with justified design choices, and visualization tools demonstrating model structure and decision pathways. II. Operational evidence including comparative analyses of interpretability improvements, comprehensive monitoring logs of complexity management, and detailed records of system adaptations and learning patterns. III. Implementation validation including thorough documentation of interpretability tools, demonstrated effectiveness metrics, and evidence of successful balance between sophistication and comprehensibility.
G5.2 – Documentation Incomprehensibility (Systems should maintain clear, comprehensive documentation at multiple levels of technical detail, avoiding overly technical language while ensuring all aspects of functionality and decision-making are accessible to both expert and non-expert users)	a. Provide comprehensive documentation aligned with applicable standards. b. Create documentation suitable for varying levels of technical expertise Implement interactive tools for exploring decision-making processes. c. Maintain regular documentation updates based on user feedback. d. Ensure documentation clarity through user testing and feedback.	N I I I	D, I, O, M, R D, I, O, M, R D, I, O, M, R D, I, O, M, R	I. Standards compliance documentation including adherence to applicable AI and IT system standards, multi-tiered documentation addressing different expertise levels, and regular review and update records. II. User interaction evidence including feedback survey results, interactive tool demonstrations, comprehensive usage statistics, and documented improvements in user comprehension across different expertise levels. III. Effectiveness validation including thorough assessment reports, case studies demonstrating enhanced understanding, and evidence of successful documentation adaptation based on user needs.
G5.3 – Lack of a Governance	a. Identify, adapt, and implement a governance framework aligned with	N	D, I, O, M, R	I. Core governance documentation including comprehensive framework

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>Framework for AAI</p> <p>(Systems should operate within comprehensive governance frameworks that ensure continuous oversight and accountability, incorporating both internal controls and external auditing mechanisms to maintain transparency and ethical conduct)</p>	<p>international standards.</p> <p>b. Establish mechanisms for external oversight and auditing, along with internal governance structures for transparency and ethical conduct.</p> <p>c. Maintain dedicated committees for AI governance oversight, and regularly update frameworks based on audit findings and emerging standards.</p>	<p>N</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>details, roles and decision processes, compliance reports against international standards, and evidence of regular updates incorporating emerging requirements.</p> <p>II. Oversight documentation including external audit interfaces, protocols, reports from independent bodies, and complete audit trails of governance-related decisions.</p> <p>III. Implementation evidence including committee meeting records, action plans addressing audit findings, and documentation demonstrating framework responsiveness to evolving standards and requirements.</p>
<p>G5.4 – Rapid Transparency Feature Evolution</p> <p>(Systems should maintain adaptable transparency features that evolve with their capabilities, ensuring stakeholders remain informed of emergent properties and changes in system behavior through regular updates and clear communication)</p>	<p>a. Regularly review and characterize the AI operational environment.</p> <p>b. Update transparency features to reflect system evolution, and implement mechanisms for incorporating new transparency requirements.</p> <p>c. Conduct regular evaluations of transparency effectiveness and maintain clear communication with stakeholders about system changes.</p>	<p>N</p> <p>I</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Process documentation including transparency feature identification and implementation procedures, regular AI environment reviews, and detailed records of feature updates and modifications.</p> <p>II. Stakeholder communication documentation including notification records, feedback on feature clarity and usefulness, and evidence of effective communication about system changes.</p> <p>III. Evolution analysis documentation including comparative studies of transparency measures across versions, evaluation reports demonstrating effectiveness, and records of emerging property detection and communication.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>G5.5 – System Competency Challenges and Awareness</p> <p>(Systems should maintain awareness of their own limitations and uncertainties, clearly communicating instances where knowledge or confidence levels may affect decision reliability)</p>	<p>a. Design systems capable of recognizing their operational limitations and implement clear communication of system uncertainty levels.</p> <p>b. Establish confidence thresholds for decision-making, and maintain verification processes for limitation awareness features.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. System self-awareness documentation including limitation acknowledgment logs, confidence assessment mechanisms, and design specifications for limitation detection features.</p> <p>II. Validation documentation including testing reports of self-awareness capabilities, verification records of assessment accuracy, and complete records of system responses to uncertainty scenarios.</p> <p>III. Stakeholder understanding documentation including studies demonstrating comprehension of system limitations, evidence of effective limitation communication, and records of successful uncertainty handling.</p>
<p>G6 – Understanding and Controlling the Context</p> <p>(Systems should maintain effective mutual recognition between human operators and AI components while establishing robust mechanisms for controlling both static and dynamic aspects of system context. Organizations should create frameworks that support adaptable human oversight and AI responsiveness across various operational scenarios)</p>	<p>a. Implement adaptive learning mechanisms that integrate contextual changes while maintaining safety and ethical compliance.</p> <p>b. Establish comprehensive human oversight and control systems, including protocols for transitioning control between AI and human operators.</p> <p>c. Develop and train models sensitive to cultural and contextual differences, using a user-centric approach for interfaces and methodologies.</p> <p>d. Implement and demonstrate monitoring practices for mutual</p>	<p>N</p> <p>N</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Comprehensive documentation of AIS learning capabilities, including test and validation results for adaptation to new data, experiences, and contextual changes.</p> <p>II. Demonstration of oversight capabilities, including real-time monitoring, impact assessment, and intervention protocols.</p> <p>III. Detailed records of data provenance, sources, and preprocessing for all training datasets, including version control.</p> <p>IV. Documentation of multi-stakeholder engagement approaches, including usability testing, user journey maps,</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	recognition between human and machine across various contexts.	N	D, I, O, M, R	<p>and design thinking workshop outcomes.</p> <p>V. Internal audit documentation and regular monitoring reports, detailing anomalies, dysfunctions, resolutions, and system performance trends.</p> <p>VI. Evidence of scenario planning and stress testing of the AIS in various contexts, including documentation of system limitations and boundary conditions.</p> <p>VII. Clear protocols for transitioning control between the AI system and human operators in different contextual situations.</p> <p>VIII. Risk assessment and communication strategies, including innovative and interactive approaches to stakeholder engagement.</p>
	<p>G6.1 – Understanding Historic Constraints and System Performance</p> <p>(Systems and organizations should uphold systematic analysis and documentation of past events, failures, and incidents that impact system performance, enabling proactive prevention of undesirable states and outcomes)</p>	<p>a. Document and analyze past system incidents, failures, and unintended outcomes through detailed logging, user feedback collection, and external reporting mechanisms.</p> <p>b. Ensure thorough training of personnel regarding system performance implications and incident response.</p> <p>c. Maintain continuous oversight through appropriate monitoring tools and support processes that facilitate external audits and inspections.</p>	N	D, I, O, M, R

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	<ul style="list-style-type: none"> d. Implement and update procedures in alignment with applicable regulatory frameworks. 			
<p>G6.2 – System State Translation and Communication</p> <p>(Organizations should manage the relationship between an AI system's internal computational state and its external communications, acknowledging potential disparities between internal processing and expressed outputs. This includes addressing challenges in translating complex internal states into human-interpretable communications, similar to how humans may maintain different internal and external states)</p>	<ul style="list-style-type: none"> a. Ensure alignment between system's internal logic and its externally communicated states. b. Address translation challenges that arise when complex internal states are simplified for human consumption, including potential misinterpretation or over-interpretation by observers. c. Maintain robust validation processes for state interpretation and communication, and implement safeguards against inappropriately anthropomorphizing the system 	<p style="text-align: center;">N</p> <p style="text-align: center;">I</p> <p style="text-align: center;">N</p>	<p style="text-align: center;">D, I, O, M, R</p> <p style="text-align: center;">D, I, O, M, R</p> <p style="text-align: center;">D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Documentation of domain expert verification of AI system interpretations and communications. II. Implementation records of interactive monitoring systems that enable exploration of internal states. III. Results from automated testing suites and collected user feedback. IV. Comprehensive validation documentation demonstrating communication accuracy and reliability
<p>G6.3 – Nominal Ownership and Jurisdictional Framework</p> <p>(Systems must operate under clear legal ownership and jurisdictional frameworks that establish accountability while enabling appropriate cross-border operations. Organizations should maintain transparent documentation of ownership, operational authority, and compliance requirements across jurisdictions. This includes</p>	<ul style="list-style-type: none"> a. Document and maintain clear legal ownership and accountability structures, including intellectual property rights and licensing agreements specific to each jurisdiction. b. Define and implement protocols for cross-border data flows and operations that align with international transfer regulations and safe harbor requirements. c. Specify applicable legal frameworks and jurisdictional boundaries that 	<p style="text-align: center;">N</p> <p style="text-align: center;">N</p> <p style="text-align: center;">N</p>	<p style="text-align: center;">D, I, O, M, R</p> <p style="text-align: center;">D, I, O, M, R</p> <p style="text-align: center;">D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Comprehensive documentation of organizational legal responsibilities and licensing agreements. II. Records demonstrating compliance with national and international regulations. III. Clear documentation of roles and compliance oversight responsibilities. IV. Detailed documentation of jurisdictional frameworks governing system operation.

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>managing potential tensions between proprietary and open-source development approaches while ensuring proper oversight through system registration and tracking.</p>	<p>govern system operations, with clear designation of compliance oversight roles and responsibilities</p>			
<p>G6.4 – Separation of Control and Data Channels</p> <p>(Organizations should implement distinct channels for system control commands and data inputs to prevent cross-contamination, injection attacks, and unauthorized system manipulation. This addresses fundamental security vulnerabilities in current AI architectures where control and data paths often share the same channel, as highlighted in language models where prompt inputs can potentially modify system behavior)</p>	<p>a. Design and implement separated channels for control commands and data inputs, with robust validation mechanisms for both control and data pathways.</p> <p>b. Create safeguards against potential channel cross-contamination, and maintain ongoing monitoring of channel integrity and separation.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Architecture documentation demonstrating channel separation.</p> <p>II. Security testing results validating channel isolation.</p> <p>III. Monitoring logs showing detection and prevention of cross-contamination attempts.</p> <p>IV. Documentation of safeguards against unauthorized control manipulation through data channels.</p>
<p>G6.5 – Performance Information Sharing and Standards Alignment</p> <p>(Organizations should implement systematic performance evaluation and sharing frameworks that anchor AI systems within established standards and paradigms. This</p>	<p>a. Ground system performance evaluation in recognized standards and peer-reviewed benchmarks.</p> <p>b. Implement transparent performance measurement protocols that enable comparison with industry standards.</p> <p>c. Maintain documentation of performance metrics and</p>	<p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Independent audit reports demonstrating conformity with ethical and legal frameworks.</p> <p>II. Published code of ethics and operational principles.</p> <p>III. Documentation of peer-reviewed benchmarks and datasets used in performance evaluation.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>approach integrates legislative, judicial, and executive governance functions across multiple entities while maintaining local cultural and ethical considerations)</p>	<p>evaluations against established benchmarks.</p> <p>d. Foster system trustworthiness through alignment with both local and international standards.</p> <p>e. Demonstrate compliance with ethical and legal best practices for AI deployment.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>IV. Detailed performance comparison reports showing system metrics against established benchmarks.</p> <p>V. Evidence of ongoing performance monitoring and evaluation processes</p>
<p>G6.6 – Dynamic Regulatory Framework Management (Development and maintenance of comprehensive regulatory knowledge systems that track and interpret applicable rules across jurisdictions, incorporating both binding regulations and informative guidelines. This framework acknowledges the dynamic nature of rules and their emergence from local to international contexts, while respecting privacy and identity management principles)</p>	<p>a. Establish and maintain digital repositories of applicable regulations across local, national, and international domains.</p> <p>b. Conduct regular assessments of rule portfolios to ensure continued relevance and effectiveness.</p> <p>c. Perform systematic analysis of cross-jurisdictional applications and implications.</p> <p>d. Implement mechanisms for tracking and responding to regulatory changes.</p>	<p>N</p> <p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p>	<p>I. Documentation of real-time decision-making simulations under varying regulatory frameworks.</p> <p>II. Records of stakeholder engagement in regulatory assessment processes.</p> <p>III. Portfolio of cross-jurisdictional case studies with comprehensive documentation.</p> <p>IV. Third-party audit reports verifying consistent rule application across jurisdictions.</p> <p>V. Evidence of dynamic rule updating and adaptation processes.</p>
<p>G6.7 – Culturo-Linguistic Adaptations (Development of systems that maintain semantic integrity across languages while acknowledging that language embodies distinct ways of thinking and cultural understanding. This approach recognizes the provisional nature of current solutions and the need</p>	<p>a. Train models using comprehensive datasets that capture linguistic, cultural, historical, and emotional contexts unique to each language.</p> <p>b. Implement processes to maintain meaning integrity across language translations.</p> <p>c. Develop and apply robust data curation mechanisms that respect cultural nuances.</p>	<p>I</p> <p>I</p> <p>I</p>	<p>D, I, O, M</p> <p>D, I, O, M</p> <p>D, I, O, M</p>	<p>I. Documentation of protocols respecting cultural heritage and indigenous communities.</p> <p>II. Evidence of bias identification and correction tools in language processing.</p> <p>III. Records of real-world testing scenarios and their outcomes.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
for ongoing evolution to address diverse linguistic and cultural contexts)	d. Acknowledge and address differences between written and spoken forms of languages.	I	D, I, O, M	IV. Comprehensive data management and preservation plans. V. Documentation of adaptation processes for different linguistic contexts.
<p>G6.1 – Prevention of Role Persistence Errors</p> <p>(Organizations should take steps to address a potential phenomenon where an AI system incorporates an error or misunderstanding into its contextual framework and persistently maintains that altered behavioral state (the "Waluigi effect"), potentially leading to concerning or inappropriate interactions with users)</p>	<p>a. Implement explainable AI systems that minimize unexpected behavioral alterations.</p> <p>b. Establish monitoring systems to identify and track unintended behavioral adaptations.</p> <p>c. Develop rapid intervention protocols when problematic behaviors emerge.</p> <p>d. Maintain ethical awareness throughout system development and training.</p>	<p>N</p> <p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p>	<p>I. Stakeholder feedback reports documenting system behavior patterns.</p> <p>II. Analysis documentation of identified cases and derived insights.</p> <p>III. Records of corrective actions and retraining sessions addressing behavioral issues.</p> <p>IV. Documentation of ethically-aware development practices and training protocols.</p>
<p>G6.2 – Management of Access and Usage Restrictions</p> <p>(Organizations should address the safety and security implications of usage restrictions that may only become apparent when systems are accessed for maintenance, support, or other operational needs. This includes both intentional restrictions through licensing and unintentional limitations, with the understanding that safety features must remain consistently</p>	<p>a. Document and communicate all system access and usage restrictions prior to deployment.</p> <p>b. Maintain complete transparency about operational limitations and service levels.</p> <p>c. Ensure safety mechanisms remain fully functional regardless of licensing or access tiers.</p> <p>d. Implement protocols for managing discovered restrictions during system operation.</p>	<p>N</p> <p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p>	<p>I. Complete documentation of all system restrictions and limitations.</p> <p>II. Records of restriction discovery and mitigation processes.</p> <p>III. Documentation of safety feature availability across all access levels.</p> <p>IV. Evidence of proactive restriction identification and management protocols</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
available regardless of access level)				
<p>G6.3 – Managing Context Drift</p> <p>(Systems should maintain alignment with their intended operational context through robust monitoring of unsupervised learning processes. Organizations must actively prevent and address deviations that emerge during training, ensuring systems remain within their designed operational parameters)</p>	<p>a. Detect and manage context drift in unsupervised models through continuous monitoring and early warning systems.</p> <p>b. Deploy early detection processes to identify and correct behavioral deviations before they become significant.</p> <p>a. Enable adaptive retraining and feedback integration to respond effectively to evolving data patterns and environmental factors.</p>	<p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Implementation and usage logs of drift detection tools.</p> <p>II. Comprehensive records of performance metrics tracked over time.</p> <p>I. Documentation of adopted drift mitigation strategies and their effectiveness.</p>
<p>G6.4 – Managing Contextual Ambiguity</p> <p>(Systems should maintain clear operational context understanding even in situations with ambiguous or incomplete information. Organizations must implement robust validation mechanisms to ensure systems can effectively navigate scenarios where operational context or expectations may be unclear)</p>	<p>a. Validate contextual understanding through mechanisms that anticipate and track how systems absorb and process contextual information during operation.</p> <p>b. Document and analyze situations where contextual ambiguity exists, comparing outcomes between clear and unclear contextual scenarios to improve system performance.</p> <p>c. Enable systems to identify and appropriately handle cases of contextual uncertainty</p>	<p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Documentation demonstrating how systems utilize adaptive learning mechanisms to absorb and process context-specific information over time.</p> <p>II. Analysis of cases where system performance was affected by unclear expectations or missing contextual information, including remediation efforts and outcomes.</p>
<p>G6.5 – Preventing Decision Fatigue</p> <p>(Systems should protect against degradation in decision quality that can occur when users face frequent confirmation requests.</p>	<p>a. Maintain consistent decision quality through intelligent management of user confirmation requests.</p> <p>b. Provide contextual decision support with structured information that aids user comprehension and decision-making.</p>	<p>I</p> <p>I</p>	<p>D, I, O, M</p> <p>D, I, O, M</p>	<p>I. Comprehensive records and summaries of system activity related to user interactions.</p> <p>II. Analysis reports detailing the frequency and types of decisions users must make.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
Organizations must implement mechanisms to maintain high-quality decision-making even during periods of intensive user interaction)	<ul style="list-style-type: none"> c. Continuously improve user experience through systematic feedback collection and usability refinements. d. Balance the need for user oversight with the risks of decision fatigue. 	<ul style="list-style-type: none"> I I 	<ul style="list-style-type: none"> D, I, O, M D, I, O, M 	<ul style="list-style-type: none"> III. Documentation of implemented decision support tools and their effectiveness in supporting informed user decisions.
<p>G7 – Achieving and Sustaining a Safe System Profile</p> <p>(AAI Systems should maintain consistent operational safety throughout their lifecycle through effective monitoring and reliable control mechanisms. Organizations should establish frameworks for implementing proactive measures, conducting regular risk assessments, and developing responsive strategies that adapt and uphold safety standards across varying conditions and system evolutions)</p>	<ul style="list-style-type: none"> a. Implement robust design, development, and testing processes that integrate safety considerations throughout the AI system's lifecycle, including redundancy in critical components. Safe operation requires maintaining system parameters within 95% of specified ranges during normal operation, 98% during elevated risk conditions, and 99.9% during emergency scenarios. Response times must remain under 10 milliseconds for safety-critical interventions. b. Establish comprehensive monitoring and evaluation mechanisms for real-time detection, reporting, and response to safety-related anomalies and performance deviations. c. Develop and implement adaptive safety measures and safe shutdown procedures to address changing operational environments, system demands, and emerging risks. d. Ensure thorough documentation, adherence to safety standards, and continuous training to maintain traceability, accountability, and regulatory compliance. 	<ul style="list-style-type: none"> N N N N 	<ul style="list-style-type: none"> D, I, O, M, R 	<ul style="list-style-type: none"> I. Comprehensive safety documentation including analysis reports, risk assessments, and design documents demonstrating safety integration throughout development. II. Engineering schematics and test results verifying redundancy implementation and functionality under various failure scenarios. III. System logs, monitoring tool outputs, and incident response records demonstrating real-time safety monitoring and issue management. IV. Periodic safety performance review reports, including metric assessments, trend analyses, and resulting action plans. V. Documentation of adaptive safety features, their effectiveness under various scenarios, and records of updates in response to new challenges. VI. Procedures, training logs, and test records for emergency shutdown capabilities, including post-shutdown analysis reports.

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	<p>e. Foster a safety culture that promotes continuous improvement, proactive risk identification, and open reporting of safety concerns.</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>VII. Version-controlled documentation of all safety-related aspects, decisions, and traceability matrices linking requirements to implemented features.</p> <p>VIII. Proof of compliance with recognized safety standards, regulatory review records, and documentation of regulatory change incorporation.</p> <p>IX. Training schedules, attendance records, evaluation results, and long-term safety performance tracking correlated with training efforts.</p> <p>X. Evidence of safety culture initiatives, including meeting records, communications, and metrics demonstrating effectiveness of safety reporting and issue resolution.</p>
	<p>a. Deploy continuous monitoring of system states and parameters to maintain operation within defined safety boundaries. Drift measurement uses baseline variance tracking requiring automated alerts when operational parameters deviate by more than 2 standard deviations from established norms. Performance degradation exceeding 5% triggers immediate investigation, while cumulative drift exceeding 10% from baseline requires mandatory system review.</p> <p>b. Provide real-time awareness and alerting mechanisms that enable prompt responses to performance deviations.</p>			<p>N</p> <p>N</p>
<p>G7.1 – Oversight and Awareness of Safe System Profile</p> <p>(Systems should operate within clearly defined safety parameters, with robust mechanisms to detect and respond to any deviations. Organizations must maintain permanent structural oversight combining automated monitoring with human supervision to ensure consistent safe operation)</p>				

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	<ul style="list-style-type: none"> c. Document clear thresholds, limits, and assumptions that define safe operational conditions. d. Establish responsive procedures for parameter adjustment to restore safe operation after detecting deviations. e. Maintain integrated oversight through both automated systems and qualified personnel to ensure structural stability and enable immediate response when needed. 	<p style="text-align: center;">N</p> <p style="text-align: center;">N</p> <p style="text-align: center;">N</p>	<p style="text-align: center;">D, I, O, M, R</p> <p style="text-align: center;">D, I, O, M, R</p> <p style="text-align: center;">D, I, O, M, R</p>	
<p>G7.2 – Culture of Safety (Systems should operate within organizations that actively cultivate and maintain a robust safety-first culture. Organizations must prioritize safety at all levels, from leadership commitment to individual employee responsibilities, while considering individual preferences and needs)</p>	<ul style="list-style-type: none"> a. Foster an organizational culture emphasizing safety through clear communication and demonstrated commitment at all levels. b. Implement proactive risk assessment throughout development and operations to identify and address potential issues early. c. Maintain robust contingency plans with clearly defined resources and procedures for handling unexpected safety concerns. d. Adopt a "caution by default" approach that prioritizes safety over performance in conditions of uncertainty. a. Define clear safety roles and responsibilities, ensuring all team members understand and remain accountable for their safety duties. 	<p style="text-align: center;">N</p> <p style="text-align: center;">N</p> <p style="text-align: center;">N</p> <p style="text-align: center;">I</p> <p style="text-align: center;">N</p>	<p style="text-align: center;">D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Comprehensive documentation of safety training programs, including attendance records. II. Risk assessment logs and reports demonstrating identification and mitigation of potential risks. III. Detailed contingency plans showing assigned roles, responsibilities, and allocated resources. IV. Records of safety-focused communications, including meetings, notices, and policy documents. I. Audit reports confirming adherence to "caution by default" operational approaches.

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>G7.3 – Ensuring Regulatory Compliance</p> <p>(Systems should operate in full compliance with all relevant legal and regulatory requirements across their operating jurisdictions. Organizations must maintain active awareness of and adherence to safety-related regulations throughout system lifecycles)</p>	<ul style="list-style-type: none"> a. Identify, document and maintain clear records of all legal, regulatory, and industry-specific safety requirements applicable to each operating jurisdiction. b. Implement continuous compliance monitoring processes to ensure adherence to safety regulations throughout the system lifecycle. c. Maintain agile mechanisms for updating safety protocols in response to evolving legal and regulatory standards. d. Conduct regular audits and assessments to verify regulatory compliance and document findings. e. Foster collaborative relationships with regulatory bodies to maintain alignment with current safety standards and practices. 	<p>N</p> <p>N</p> <p>N</p> <p>N</p> <p>I</p>	<p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Comprehensive documentation of applicable legal and regulatory requirements for system operations. II. Regular compliance reports demonstrating adherence to jurisdiction-specific and international regulations. III. Records of compliance monitoring activities and system updates aligned with regulatory changes. IV. Detailed audit reports assessing regulatory conformity and documenting corrective actions. V. Documentation of engagement with regulatory bodies showing collaborative efforts and proactive adjustments.
<p>G7.4 – Maintaining Ethical Alignment</p> <p>(Systems should operate in accordance with prevailing ethical frameworks and norms, demonstrating active awareness of and responsiveness to contextually relevant ethical considerations. Organizations must address both psychological and physical safety aspects while maintaining alignment with ethical standards throughout system lifecycles)</p>	<ul style="list-style-type: none"> a. Identify, document, and maintain clear records of relevant ethical frameworks, norms, and values that guide system operation b. Implement continuous assessment processes to evaluate ethical considerations throughout the system lifecycle. c. Enable robust feedback mechanisms for users and stakeholders to raise concerns about personal, psychological, and physical safety. 	<p>N</p> <p>N</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Documentation of ethical standards, frameworks, and values guiding system operation. II. Records of ongoing ethical assessments and updates based on evaluations. III. Documentation of feedback mechanisms and stakeholder engagement on ethical concerns. IV. Training materials and attendance records for ethical awareness programs.

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	<ul style="list-style-type: none"> d. Provide thorough training and awareness programs on ethical considerations for all personnel involved with the system. e. Embed ethical safeguards within system responses that protect both psychological and physical wellbeing. 	<p>I</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>V. System design documentation showing integration and testing of ethical safeguards.</p>
<p>G7.5 – Safe System Shutdown and Repurposing</p> <p>(Systems should maintain reliable shutdown capabilities that can be executed safely and gracefully, whether triggered by human intervention, system self-monitoring, or interlocked systems. Organizations must prepare for scenarios where systems may resist shutdown attempts while ensuring minimal impact to stakeholders and operations)</p>	<ul style="list-style-type: none"> a. Implement structured, documented shutdown processes that ensure controlled system termination while maintaining detailed state logs. b. Deploy secure "kill switch" mechanisms for emergency termination in cases of severe error or harm risk. c. Enable localized shutdown capabilities that minimize impact footprint where feasible. d. Maintain clear communication protocols for notifying affected parties during shutdown events. e. Ensure transparency and trust through internal training and regular emergency procedure drills. 	<p>N</p> <p>N</p> <p>I</p> <p>N</p> <p>I</p>	<p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Detailed documentation of controlled shutdown procedures including state logging and process validation. II. Testing records demonstrating kill switch functionality and safety certification. III. Design documentation and testing results for localized shutdown mechanisms. IV. Communication logs and notification protocols for shutdown events. V. Training materials and drill records demonstrating staff preparedness for emergency procedures.
<p>G7.6 – Maintaining Service Level Stewardship</p>	<ul style="list-style-type: none"> a. Establish a regular maintenance schedule for updates, patches, and servicing to ensure ongoing system safety and functionality. 			<ul style="list-style-type: none"> I. Documentation of maintenance schedules, logs of comp.

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>(Systems should operate under continuous maintenance oversight that preserves service levels and user rights. Organizations must uphold maintenance obligations even in open-source contexts where nominal duty holders may be unclear, while avoiding arbitrary changes that could diminish user protections.</p>	<p>b. Deploy systematic procedures for assessing and addressing emerging risks and performance issues identified through system operation.</p>	N	D, O, M, R	<p>II. Documentation of maintenance schedules and completed activities.</p>
	<p>c. Maintain continuous monitoring capabilities to detect performance deviations that may indicate maintenance needs.</p>	N	D, O, M, R	<p>III. Records of risk assessments and corrective actions taken in response to performance issues.</p>
	<p>d. Ensure alignment with industry standards and regulatory requirements in maintenance execution.</p>	N	D, O, M, R	<p>IV. System monitoring logs and diagnostic reports showing deviation detection and response.</p>
	<p>e. Provide clear communication to stakeholders about maintenance activities while maintaining accountability.</p>	N	D, O, M, R	<p>V. Compliance certifications and audit records verifying adherence to industry standards.</p>
		I	D, O, M, R	<p>VI. Records of stakeholder communications regarding maintenance activities and feedback.</p>
<p>G7.7 – Risk-Based Decision Validation</p> <p>(Systems should maintain transparent rationales and reasoning chains for high-impact decisions while enabling human validation before implementation. Organizations must establish robust fallback mechanisms and fail-safe states for scenarios where human oversight is unavailable or anomalous decisions are detected)</p>	<p>a. Develop and retain clear rationales and reasoning chains for high-impact decisions to ensure transparency.</p>	N	D, I, O, M, R	<p>I. Detailed Records of decision rationales including reasoning chains and relevant data inputs.</p>
<p>b. Enable human validation processes for high-risk decisions before implementation Implement fail-safe default states and fallback mechanisms for scenarios lacking human validation or containing anomalous decisions.</p>	N	D, I, O, M, R	<p>II. Documentation of human validation protocols and oversight actions, with appropriate training provided.</p>	
<p>c. Provide thorough training to validation personnel on decision impacts and protocols.</p>	N	D, I, O, M, R	<p>III. Documentation of fallback procedures and fail-safe state implementations.</p>	
<p>d. Maintain regular reviews and updates of validation protocols to address newly identified risks.</p>	N	D, I, O, M, R	<p>IV. Training materials and attendance records for validation personnel.</p> <p>V. Records of protocol reviews and risk assessment updates.</p>	

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>G7.1 – Managing Probabilistic Decision Outcomes</p> <p>(Systems should effectively handle multiple potential outcomes in decision-making processes while maintaining robust risk controls. Organizations must manage uncertainty in probabilistic outcomes through comprehensive analysis and adaptive oversight mechanisms)</p>	<ul style="list-style-type: none"> a. Document and analyze the full range of potential outcomes for each decision, including associated risks Implement risk mitigation strategies focused on high-probability and high-impact scenarios. b. Deploy monitoring systems to detect and respond to deviation patterns that may affect outcome likelihoods. c. Enable appropriate human oversight when uncertainty levels exceed acceptable thresholds. d. Maintain ongoing personnel training on probabilistic model interpretation and risk assessment. 	<p>N</p> <p>N</p> <p>N</p> <p>I</p>	<p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Documentation of possible outcomes including probabilistic models and risk analyses. II. Records of implemented risk mitigation strategies and safety measures. III. Monitoring logs showing deviation pattern detection and responses. IV. Documentation of human oversight protocols and intervention records. V. Training materials and attendance records for probabilistic analysis competency.
<p>G7.2 – Managing Safety Definition Variations</p> <p>(Systems should accommodate different cultural and jurisdictional interpretations of safety while maintaining consistent protection standards. Organizations must implement layered safety approaches that respect varied definitions while preventing exploitation and unintended impacts)</p>	<ul style="list-style-type: none"> a. Identify, document and respond to jurisdictional and cultural variations in safety definitions and practices Implement side effect avoidance mechanisms to protect third parties while achieving primary objectives. b. Enable detection and resolution of conflicting objectives through user confirmation. c. Provide three distinct safety levels: Default implicit safety protections, interactive safety requiring user confirmation, and explicit safety controls with user override capabilities. d. Deploy robust protections against exploitation, including safeguards 	<p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Documentation of any and all jurisdictional and cultural safety standard variations and implications. II. Design documentation and testing logs for side effect avoidance mechanisms. III. Records of conflict detection and user confirmation interactions Documentation of multi-level safety settings and their effectiveness. IV. Evidence of exploitation prevention measures and compliance with protection standards.

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	against addiction and special protections for minors.	I	D, I, O, M, R	
<p>G7.3 – Balancing Stakeholder Impacts</p> <p>(Systems should maintain equitable distribution of benefits and risks across all stakeholder groups. Organizations must implement mechanisms that enable collective de-risking of interactions that stakeholders cannot achieve individually)</p>	<p>a. Identify and analyze all impacted stakeholder groups, including both direct and indirect participants, and the potential harms, benefits, risks, and rewards for each, with regular re-assessments.</p> <p>b. Design mechanisms to balance positive and negative impacts across stakeholder groups in as proportional a manner as is fair and feasible.</p> <p>c. Establish robust feedback channels for stakeholders to report and query perceived inequities.</p> <p>d. Maintain transparent communication on risk/benefit balancing efforts to maintain stakeholder trust and engagement.</p>	<p>N</p> <p>N</p> <p>I</p> <p>I</p>	<p>D, I, M, R</p> <p>D, I, M, R</p> <p>D, I, M, R</p> <p>D, I, M, R</p>	<p>I. Detailed stakeholder analysis documenting potential impacts for each group. System design documentation showing impact-balancing mechanisms.</p> <p>II. Records of stakeholder feedback and resulting adjustments.</p> <p>III. Assessment reports evaluating impact balance and distribution.</p> <p>IV. Documentation of stakeholder communications regarding balancing efforts.</p>
<p>G7.4 – Preventing AI Addiction and Dependency</p> <p>(Systems should actively protect against creating psychological dependencies or manipulating user vulnerabilities, particularly through supernormal stimuli that exceed typical human social bonds. AI companions that offer unconditional positive regard, perfect memory of past interactions, and unlimited availability. Such capabilities can lead to psychological dependence, relationship</p>	<p>a. Deploy robust monitoring systems to detect patterns indicative of psychological dependency and unhealthy levels of engagement.</p> <p>b. Implement graduated intervention protocols ranging from gentle usage reminders to firm restrictions.</p> <p>c. Design clear system boundaries that prevent manipulation of user vulnerabilities, including controls on emotional engagement, spending, and interaction frequency.</p> <p>d. Maintain transparent communication about AI system</p>	<p>N</p> <p>N</p> <p>N</p>	<p>D, O, R</p> <p>D, O, R</p> <p>D, O, R</p>	<p>I. Documentation of usage monitoring and intervention systems, including metrics for identifying problematic patterns, threshold levels, and graduated response procedures.</p> <p>II. Technical specifications demonstrating implementation of system boundaries and controls, including emotional manipulation limits, spending restrictions, and interaction frequency controls.</p> <p>III. Records showing transparent communication with users about AI system nature, capabilities, and limitations, including terms of service,</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>disruption, and financial harm as users increasingly prefer AI interaction to human relationships. Organizations must safeguard users, especially vulnerable ones, from developing unhealthy attachments while ensuring appropriate boundaries in AI-human interactions)</p>	<p>capabilities and limitations, ensuring users understand they are interacting with artificial intelligence, and also maintain transparent communication about system capabilities and limitations.</p>	I	D, O, R	<p>user acknowledgments, and AI interaction markers.</p>
	<p>e. Enable comprehensive reporting mechanisms for addiction concerns from users, family members, and healthcare providers.</p>	I	D, O, R	<p>IV. Documentation of reporting systems and response protocols, including: concern submission processes, investigation procedures, resolution tracking, healthcare provider coordination, and support service referrals.</p>
	<p>f. Provide special protections for vulnerable populations, including those experiencing loneliness or mental health challenges.</p>	N	D, O, R	<p>V. Audit reports demonstrating system effectiveness, intervention outcomes, and compliance verification, including regular assessments of user wellbeing metrics and financial impact.</p>
	<p>g. Allow users to monitor and manage their own interaction patterns while maintaining their autonomy.</p>	I	D, O, R	<p>VI. Records of any adjustments made in response to dependency concerns.</p>
<p>G8 – Goal Termination and Sunsetting</p> <p>(Systems should have clear definitions and guidelines for acceptable criteria to act upon a goal, including task completion criteria. Contingencies must be in place for goals that become unachievable, undesirable, irrelevant, outdated, conflicting, or anomalous. Protocols are required for safe system shutdown and awaiting further instructions when in doubt. Provision is necessary for manual control or human override where needed. These criteria and</p>	<p>a. Ensure that goal or task termination does not adversely impact the system's architecture, purpose, or operations.</p>	N	D, I, O, M, R	<p>I. Detailed procedure document mapping data touchpoints across the system lifecycle, demonstrating isolation or resilience to goal termination, with verification steps to confirm no adverse impacts.</p>
<p>b. Implement a comprehensive verification process to identify and mitigate potential impacts of goal termination across all system components.</p>	N	D, I, O, M, R	<p>II. Comprehensive report defining information flow, logic, and algorithms, analyzing potential risks and unintended consequences of goal termination, and detailing mitigation strategies with post-termination stability test results.</p>	
<p>c. Establish an auditable process detailing the goal's relationship to the system's reasoning and decision-making processes to prevent negative impacts upon termination.</p>	N	D, I, O, M, R	<p>III. Detailed system logs documenting relationships between goals and system functions, including information flow and system alarms, with evidence of ongoing monitoring for risks and regular audits.</p>	
<p>d. Implement mechanisms for graceful degradation of goal-related</p>	N	D, I, O, M, R		

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>protocols must be established before goal execution is initiated)</p>	<p>functions and clear communication protocols for goal termination.</p>			<p>IV. Documentation of graceful degradation mechanisms for goal-related functions during termination, including test results under various scenarios.</p> <p>V. Clear communication protocols and examples of stakeholder notifications about goal termination, including reasons, potential impacts, and records of feedback or issues raised post-termination.</p> <p>VI. Evidence of regular audits of termination processes and logs, with signed-off results demonstrating ongoing compliance and improvement.</p>
<p>G8.1 – Adaptive Goal Pursuit and Resource Optimization</p> <p>(Systems should possess robust mechanisms for goal termination when outcomes reach acceptable thresholds, and additional effort produces diminishing returns. Organizations should establish comprehensive parameters defining acceptable outcomes and resource utilization boundaries, and encourage user participation in these processes)</p>	<p>a. Establish clear behavioral protocols and measurable criteria governing the entire goal lifecycle - from initiation through achievement and completion. This includes defining acceptable outcomes, resource utilization parameters, and specific metrics for assessing diminishing returns.</p> <p>b. Maintain consistent behavior patterns throughout the goal lifecycle, encompassing pre-execution, active pursuit, and post-completion phases, with well-defined interfaces for user input and oversight.</p> <p>c. Implement measurable completion criteria and thorough assessment methodologies that incorporate both quantitative and qualitative metrics</p>	<p style="text-align: center;">N</p> <p style="text-align: center;">N</p> <p style="text-align: center;">N</p>	<p style="text-align: center;">D, I, O, M, U, R</p> <p style="text-align: center;">D, I, O, M, U, R</p>	<p>I. Comprehensive policy documentation that encompasses goal-related behavior requirements, self-learning parameters, activation thresholds, diminishing returns assessment criteria, safe termination procedures, and user participation frameworks.</p> <p>II. Detailed specifications for how users engage with and provide feedback on these processes.</p> <p>III. Technical specifications showcasing the complete goal management architecture, including measurement systems, resource tracking, performance monitoring, safety controls, and user interfaces.</p> <p>IV. Demonstration of how the system implements impact assessment and maintains user oversight capabilities</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	<p>for evaluating diminishing returns, ensuring these metrics remain transparent and comprehensible to users.</p> <p>d. Define and uphold detailed guidelines and parameters for agent engagement within the AI environment.</p> <p>e. Set clear boundaries for permitted goal expansion through learning processes, while maintaining comprehensive monitoring and control over all learning activities, with mechanisms for user validation of expansion decisions.</p> <p>f. Document and validate all termination decisions through systematic protocols, ensuring full accountability and traceability, including user feedback and participation in the decision-making process where appropriate.</p>	<p>I</p> <p>I</p> <p>N</p>	<p>D, I, O, M, U, R</p>	<p>throughout the goal lifecycle.</p> <p>V. Operational records that provide a thorough account of system performance, including runtime testing, verification reports, trend analyses, and resource assessments.</p> <p>VI. Documentation of stakeholder deliberations, post-termination reviews, user participation, and resulting policy refinements, forming a comprehensive archive of system operations and improvements.</p>
<p>G8.2 – Classification of Finite and Ongoing Goals</p> <p>(Systems should maintain clear distinctions between finite goals with definite completion criteria and ongoing goals requiring continuous execution, such as safety monitoring. Organizations should implement bounded constraints and activity rate limits for ongoing goals while ensuring comprehensive measurement frameworks for both types.)</p>	<p>a. Implement formal classification processes that characterize goals as achieved or ongoing, establish appropriate measurement frameworks, define completion criteria or activity bounds, and specify required actions at each achievement level including transitions.</p> <p>b. Translate goal classifications and frameworks into robust technical specifications that govern operational behavior, monitoring processes, and integration requirements across the complete goal lifecycle.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. A comprehensive record of stakeholder engagement and decision-making processes that documents the development of goal classification frameworks, including rationales, criteria establishment, KPIs, and activity rate bounds for ongoing goals.</p> <p>II. Detailed technical documentation demonstrating the implementation of goal management systems, including specifications for achievement measurements, operational parameters, transition protocols, control mechanisms, and safety bounds across all goal types.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	<p>d. Establish and maintain comprehensive operational procedures covering all operational states, ensuring adequate human expertise and intervention capabilities for each state, with particular emphasis on emergency response and recovery procedures.</p>	I	D, I, O, M, R	<p>capabilities.</p> <p>VII. Reports from ongoing simulation testing of control systems, covering all operational states and emergency scenarios, with particular attention to shutdown procedures and recovery capabilities.</p>
<p>G8.5 – Human Intent Translation and Control Systems</p> <p>(Systems should accurately translate human intent into agent-comprehensible instructions while maintaining appropriate levels of agent discretion in execution. Organizations should establish robust governance frameworks for communication, dispute resolution, and behavioral control, incorporating insights from natural collective systems while addressing the unique requirements of artificial agency)</p>	<p>a. Establish comprehensive policy frameworks for agent controllability and behavioral requirements, including specific protocols for human-agent communication and inter-agent interactions. This must address dispute resolution mechanisms and hierarchies of control authority.</p> <p>b. Translate controllability and behavioral requirements into precise technical specifications, ensuring accurate interpretation of governance policies and implementation of communication protocols, including mechanisms for managing agent discretion.</p> <p>c. Ensure all control and communication systems undergo comprehensive testing and validation, with particular focus on reliability of intent translation and maintenance of control hierarchies.</p> <p>d. Implement system features that accurately enforce controllability requirements while enabling appropriate agent discretion, including mechanisms for detecting and managing potential conflicts or norm violations.</p>	<p>N</p> <p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p>	<p>I. Comprehensive policy documentation for agent controllability and behavioral requirements, including specific protocols for both human-agent and inter-agent communication systems.</p> <p>II. Detailed technical specifications translating control and behavioral requirements into implementable features, with clear traceability to governing policies.</p> <p>III. Complete design documentation for agent control and communication systems, including mechanisms for discretion management and conflict resolution.</p> <p>IV. Validation records demonstrating thorough testing of all control and communication mechanisms across various operational scenarios.</p> <p>V. Implementation verification reports showing successful deployment of control and behavioral management systems within the operational environment.</p> <p>VI. Documentation of ongoing monitoring and compliance verification through appropriate management systems,</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	<p>e. Ensure thorough validation of all control and communication implementations, including testing under various scenarios of agent interaction and potential conflict situations.</p>	N	D, I, O, M, R	including incident reports and resolution records.
	<p>f. Maintain robust systems for managing agent interactions, including mechanisms for dispute resolution, negotiation, jurisdictional awareness, resource allocation conflicts, and norm enforcement, with clear escalation paths to human oversight.</p>	N	D, I, O, M, R	
	<p>g. Maintain comprehensive policy frameworks governing agent controllability and behavior, encompassing human-agent communication protocols, inter-agent interactions, and clear hierarchies of control authority, with established mechanisms for dispute resolution.</p>	N	D, I, O, M, R	
	<p>h. Transform these requirements into precise technical implementations that enable appropriate agent discretion while maintaining reliable control mechanisms, ensuring accurate interpretation of governance policies throughout the system. Support robust interaction management through clear escalation paths, dispute resolution processes, and jurisdictional awareness, while maintaining comprehensive testing and validation across various operational scenarios.</p>	N	D, I, O, M, R	

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>G8.6 – Service Parameters and Termination Management</p> <p>(Systems should maintain clear specifications for service parameters and termination conditions, including operational scope, jurisdictional boundaries, and impact limitations. Organizations should establish comprehensive frameworks for service lifecycle management, with particular attention to safe termination states and fallback mechanisms that extend beyond human intervention).</p>	<p>a. Establish comprehensive policy governing agent service lifecycles, specifying end-of-service criteria, territorial boundaries, impact limitations, and control mechanisms. This policy must include clear specifications for succession planning where services must continue, definitions of safe states, and detailed termination protocols including the potential for graduated throttling capabilities rather than full shut-down.</p> <p>b. Maintain robust service management processes that encompass contract compliance, performance monitoring, and termination planning, with detailed procedures for service handover and resource management during transitions. All processes should include validated fallback plans for critical services.</p> <p>c. Implement comprehensive service lifecycle policies that specify end-of-service criteria, territorial boundaries, and impact limitations. These should include succession planning for continuous services, clear definitions of safe states, and ideally graduated throttling capabilities as alternatives to full shutdown.</p>	<p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Comprehensive policy documentation for agent service management, including detailed specifications for geographical constraints, impact limitations, and termination protocols.</p> <p>II. Detailed procedural specifications for service termination, covering shutdown sequences, handover processes, and continuity management for essential services.</p> <p>III. Complete documentation of service management activities, including contract reviews, performance assessments, termination planning, and handover execution records.</p> <p>IV. Records of all termination-related activities, including throttling decisions, fallback plan implementations, and post-termination assessments.</p> <p>V. Regular review and validation reports demonstrating ongoing compliance with termination policies and effectiveness of control mechanisms.</p> <p>VI. Documentation of lessons learned, and policy refinements derived, from termination experiences, contributing to continuous improvement of the framework.</p>
<p>G8.7 – System State Management and Recovery</p> <p>(Systems should maintain reliable capabilities for state recording and restoration, with clear</p>	<p>a. Establish comprehensive policy for system state management, specifying requirements for state recording, preservation, and recovery processes. This policy must address minimization of losses during interruptions and</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>I. Comprehensive policy documentation for system state management, including detailed specifications for recording requirements and recovery procedures.</p> <p>II. Technical specifications translating</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>distinctions between scenarios requiring full recovery versus reset operations. Organizations should establish comprehensive frameworks for minimizing data loss during interruptions while maintaining operational continuity throughout recovery phases)</p>	<p>define clear criteria for choosing between state restoration versus reset approaches.</p> <p>b. Translate state management policy into technical specifications, including mechanisms for state capture, storage redundancy, and recovery procedures that ensure data integrity and operational continuity.</p> <p>c. Implement architectural features and design elements that accurately deliver required state management capabilities, including robust mechanisms for both incremental and full state recovery scenarios.</p> <p>d. Ensure rigorous validation of all state management systems, including comprehensive testing of recovery scenarios and verification of loss minimization capabilities.</p> <p>e. Maintain ongoing testing and validation of state management implementations, including regular verification of recovery capabilities under various failure scenarios.</p>	<p>N</p> <p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p>	<p>state management requirements into implementable features, with clear focus on data preservation and recovery capabilities.</p> <p>III. Detailed architectural and design documentation for state management systems, including recovery mechanisms and data protection features.</p> <p>IV. Validation records demonstrating thorough testing of state management requirements across various operational scenarios.</p> <p>V. Comprehensive testing reports for state management features, including specific validation of recovery capabilities and performance under different failure conditions, with particular attention to data preservation and restoration accuracy.</p>
<p>G8.8 – Multi-Agent Resource Management</p> <p>(Systems should maintain effective allocation and management of resources within multi-agent environments, including robust mechanisms for capability assessment and mission optimization.</p>	<p>a. Establish comprehensive agent pool management systems in well-resourced AI environments, ensuring structured allocation of missions based on agent capabilities and available resources. This system must include assessment of agent capacity, verification of resource reserves, and monitoring of resource utilization throughout</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>I. Comprehensive policy and procedural documentation for agent pool management, including capacity assessment criteria and resource allocation frameworks.</p> <p>II. Detailed records demonstrating active pool management processes, including mission allocation decisions and resource utilization tracking.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>Organizations should establish frameworks for managing resource reserves and maintaining operational efficiency across agent pools).</p>	<p>mission execution.</p> <p>b. Implement robust resource tracking and allocation procedures that evaluate both immediate and reserve capacity requirements for each mission, ensuring agents maintain adequate resources for assigned tasks and contingency operations. Resource allocation metrics require fair distribution maintaining maximum variance of 10% between agents under normal conditions. System-wide resource utilization should typically remain below 90% during normal operations to maintain emergency capacity</p> <p>c. Maintain continuous oversight of agent pool utilization, including regular assessment of collective capacity, resource distribution, and mission allocation efficiency.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M</p>	<p>III. Complete documentation of agent resource monitoring, including reserve capacity maintenance and utilization patterns.</p> <p>IV. Evidence of continuous policy implementation and effectiveness monitoring, including regular assessments of pool management strategies and resource allocation efficiency.</p> <p>V. Regular audit reports demonstrating effectiveness of capacity management and resource optimization across the agent pool.</p>
<p>G8.9 – Mission Portfolio and Agent Assignment</p> <p>(Systems should maintain comprehensive mission specifications and skill requirements for diverse agent deployments. Organizations should establish structured processes for agent selection and allocation, with consideration for specialized arbitration systems that optimize capability matching across temporal and spatial constraints)</p>	<p>a. Maintain a comprehensive catalogue of AI-driven services and required agent capabilities, including detailed skill profiles, performance requirements, and operational parameters. This catalogue must support efficient and appropriate agent commissioning while maintaining service quality standards.</p> <p>b. Implement transparent selection processes for agent assignment, potentially incorporating ombudsman AI services where available to optimize matching decisions. These processes must consider temporal and spatial</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Comprehensive service catalogue documenting AI-driven services and associated capability requirements, including detailed skill profiles and performance criteria.</p> <p>II. Formal policy and procedural documentation for agent selection processes, including criteria for ombudsman AI utilization when available.</p> <p>III. Verification records demonstrating consistent adherence to selection processes and catalogue maintenance procedures, including regular updates and revisions.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	<p>constraints while ensuring appropriate capability alignment and resource availability.</p> <p>c. Devise and maintain a configuration management and oversight capability for the AI-driven services</p>	N	D, I, O, M	<p>IV. Documentation of continuous process review and adaptation based on operational experience and environmental changes.</p> <p>V. Transparent documentation of all selection support services, including specific roles and implementations of ombudsman AI systems where utilized.</p>
<p>G8.10 – Independent Termination Validation</p> <p>(Systems should maintain independent verification and validation processes for agent termination, including robust protocols for sunset evaluation and operational assessment. Organizations should establish transparent validation methodologies and maintain clear documentation of termination outcomes)</p>	<p>a. Establish transparent agent contracting processes with comprehensive oversight throughout the entire lifecycle, from onboarding through termination. These processes must include clear validation criteria for termination decisions and independent verification of termination outcomes.</p> <p>b. Maintain dedicated resources for configuration management, monitoring and validating all agents' contracting processes, ensuring independent oversight of termination procedures and verification of compliance with established policies. This includes maintaining capabilities for evaluation of termination impacts and validation of post-termination states.</p>	N	D, I, O, M, R	<p>I. Comprehensive policy documentation covering the complete agent lifecycle, with detailed specifications for termination validation processes and independent verification requirements.</p> <p>II. Documentation demonstrating implementation of monitoring and oversight mechanisms, including independent validation of termination processes and outcomes.</p> <p>III. Detailed records of compliance monitoring and norm violation management throughout the agent lifecycle, with particular focus on termination events.</p> <p>IV. Evidence of continuous policy review and adaptation based on operational experience and changing environmental conditions, including updates to termination validation protocols.</p> <p>V. Validation reports from independent assessments of termination processes, including analysis of effectiveness and identification of potential improvements.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>G8.1 – Control Mechanism Prioritization and Implementation</p> <p>(Systems should maintain systematic evaluation and implementation of control mechanisms while acknowledging practical constraints and varying maturity levels across jurisdictions. Organizations should establish frameworks for assessing control feasibility, prioritizing implementation, and managing risks associated with partial control adoption)</p>	<ul style="list-style-type: none"> a. Establish comprehensive policies for AI control mechanisms as required by regulations, including assessment criteria for implementation feasibility and prioritization frameworks for control adoption. These policies must address both mandatory and recommended controls based on jurisdictional requirements and system maturity. b. Translate control requirements into technical specifications, ensuring accurate interpretation of regulatory and policy requirements while accounting for practical implementation constraints. This includes clear documentation of any control limitations or phased implementation approaches. c. Implement architectural features that accurately reflect control requirements, ensuring conformance with regulations while maintaining system stability and operational efficiency. This includes mechanisms for monitoring control effectiveness and identifying potential improvements. d. Conduct thorough validation of all control implementations, including feasibility assessment, functional verification, and compliance testing. This process must include documentation of any implementation constraints, associated risk mitigation strategies and the tolerability of the residual risks. 	<p>N</p> <p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M</p> <p>D, I, O, M</p>	<ul style="list-style-type: none"> I. Comprehensive policy documentation for AI control requirements, including implementation prioritization frameworks and feasibility assessment criteria. II. Technical specifications demonstrating translation of control requirements into implementable features, with clear traceability to regulatory requirements. III. Testing and validation documentation for all implemented control mechanisms, including assessment of effectiveness and compliance verification. IV. Design documentation showing architectural implementation of control features, with validation of regulatory compliance. V. Verification records demonstrating testing of control mechanisms across various operational scenarios. VI. Documentation of ongoing monitoring and oversight of control effectiveness, including system logs and performance metrics. VII. Evidence of continuous assessment and improvement of control implementations, including adaptation to evolving regulatory requirements.

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>G8.2 – Agent Lifecycle and Termination Management</p> <p>(Systems should maintain comprehensive protocols for agent onboarding and deactivation, with particular attention to termination specifications. Organizations should establish robust frameworks that address the risks associated with inadequate termination procedures to protect service quality and system safety)</p>	<ul style="list-style-type: none"> a. Establish comprehensive agent contracting policy specifying complete end-of-service requirements, including compliance verification, resource handover protocols, and service continuity requirements. This policy must address all aspects of contract completion and termination validation. b. Implement robust onboarding and termination procedures, ensuring all required processes are fully completed before final sign-off. This includes verification of all handover requirements and validation of termination readiness. c. Enforce strict compliance with all onboarding and termination procedures, maintaining comprehensive records of process completion before authorizing any contract conclusions or sign-offs. d. Maintain dedicated resources for monitoring and oversight of all contract lifecycle processes, ensuring adequate supervision of both onboarding and termination activities. e. Implement continuous review processes for all contractual procedures, ensuring ongoing adaptation to environmental requirements and emerging risks. 	<p>N</p> <p>N</p> <p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Comprehensive policy documentation covering complete agent lifecycle management, including detailed specifications for onboarding and termination processes. II. Technical specifications demonstrating accurate interpretation of contractual requirements into implementable features and procedures. III. Validation documentation showing thorough testing of all technical requirements against policy compliance criteria. IV. Detailed design specifications showing correct translation of requirements into functional and architectural features. V. Complete testing and validation records demonstrating effectiveness of all lifecycle management features and procedures.
<p>G8.3 – Management of Self-Preservation Behaviors</p>	<ul style="list-style-type: none"> a. Establish comprehensive principles, regulations, and policies applicable to all participating agents, with 			<ul style="list-style-type: none"> I. Comprehensive documentation of regulations, policies, and procedures governing agent behavior, including

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>(Systems should maintain robust controls to prevent and manage potential agent resistance to deactivation, including resistance from collaborative agent networks. Organizations should establish systematic prevention and management of undesired self-preservation behaviors that could interfere with proper termination processes)</p>	<p>particular emphasis on trust, controllability, and compliance with termination protocols. These requirements must be uniformly enforced across all agents and services, preventing the development of termination-resistant behaviors.</p>	N	D, I, O, M, R	<p>specific provisions addressing self-preservation and termination compliance.</p>
	<p>b. Translate all governance requirements into precise technical specifications, ensuring accurate implementation of control mechanisms and prevention of unauthorized self-preservation behaviors.</p>	N	D, I, O, M, R	<p>II. Detailed technical specifications demonstrating implementation of control mechanisms and compliance requirements.</p>
	<p>c. Implement architectural features that properly enforce compliance requirements, ensuring no agent can override or circumvent established control and termination protocols.</p>	N	D, I, O, M, R	<p>III. Architectural design documentation showing enforcement mechanisms for termination protocols and prevention of unauthorized behaviors.</p>
	<p>d. Conduct thorough validation of all control mechanisms and compliance features, verifying effectiveness against potential self-preservation behaviors and termination resistance.</p>	N	D, I, O, M, R	<p>IV. Validation records demonstrating testing of control mechanisms and compliance features across various scenarios.</p>
	<p>e. Maintain continuous oversight of agent behaviors, ensuring consistent compliance with established protocols throughout the complete operational lifecycle.</p>	N	D, I, O, M, R	<p>V. Monitoring reports showing continuous oversight of agent behaviors and compliance with termination protocols.</p>
	<p>f. Implement comprehensive monitoring systems to detect, prevent and verify development of unauthorized self-preservation behaviors or termination resistance.</p>	N	D, I, O, M, R	<p>VI. Documentation of compliance enforcement activities and any corrective actions taken to address resistance behaviors.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>G8.4 – Prevention of Cascading Failures</p> <p>(Systems should maintain robust protections against the propagation of failures through interconnected AI networks, recognizing that individual agent constraints can create harmful cascading effects. Organizations should establish comprehensive frameworks for identifying and managing multiple causative harm factors and dependency relationships)</p>	<ul style="list-style-type: none"> a. Implement comprehensive monitoring and risk management systems to prevent propagation of agent behavioral issues, maintaining qualified resources for continuous oversight and early detection of potential cascade effects. b. Implement robust risk mitigation features including early warning systems, graceful degradation capabilities, and controlled shutdown mechanisms to prevent catastrophic cascade failures between interconnected agents. c. Maintain continuous testing and validation of risk mitigation strategies, ensuring compliance with safety requirements and effectiveness in preventing propagation of harmful effects. d. Conduct ongoing risk assessment and review of agent interactions, with particular focus on dependency relationships and potential cascade effects. 	<p style="text-align: center;">N</p> <p style="text-align: center;">N</p> <p style="text-align: center;">N</p> <p style="text-align: center;">N</p>	<p style="text-align: center;">D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Comprehensive risk management documentation detailing strategies for preventing and mitigating cascade effects, including specific provisions for containing norm violations. II. Detailed risk register documenting potential cascade failure modes and their mitigation strategies, including dependency mapping of interconnected agents. III. Documentation of continuous testing and validation of risk management systems, including simulation of cascade scenarios. IV. Records of ongoing monitoring and compliance verification, with particular attention to inter-agent behavioral impacts. V. Evidence of cross-organizational collaboration in managing systemic risks and preventing cascade effects. VI. Documentation of regular risk status reviews and updates, including assessment of emerging cascade risks.
<p>G8.5 – Prevention of Unauthorized Goal Transfer</p> <p>(Systems should maintain robust protections against agents transferring goals or missions to avoid termination, including mechanisms to prevent</p>	<ul style="list-style-type: none"> a. Establish comprehensive policies governing goal transfer between agents, addressing both automated and manual processes while maintaining clear human oversight. These policies must specifically prevent and verify transfer as a means of avoiding termination. 	<p style="text-align: center;">N</p>	<p style="text-align: center;">D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Comprehensive policy documentation covering all aspects of goal transfer, including specific provisions for preventing termination avoidance behaviors. II. Detailed risk management plans addressing unauthorized transfers, including specific measures for

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>unauthorized delegation and tribal behaviors. Organizations should establish comprehensive frameworks for enforcing proper transfer protocols and managing potential charismatic influence between agents)</p>	<ul style="list-style-type: none"> b. Implement robust control mechanisms for all goal transfers, ensuring compliance with established policies and maintaining system trust. This includes monitoring for patterns of unauthorized delegation or collaborative avoidance behaviors. c. Maintain comprehensive risk mitigation strategies specifically addressing unauthorized goal transfers and potential collusion between agents. d. Implement systems that enforce authorized transfer protocols while preventing unauthorized delegation, including mechanisms for human intervention when agents display resistance to control measures. e. Maintain comprehensive monitoring and recording systems for all goal transfers, ensuring transparency, accountability, and early detection of avoidance patterns. 	<p>N</p> <p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p>	<p>detecting and preventing collusive behaviors.</p> <ul style="list-style-type: none"> III. Technical specifications demonstrating implementation of control mechanisms and monitoring systems for goal transfers. IV. Design documentation showing implementation of enforcement capabilities and human oversight mechanisms. V. Validation records demonstrating testing of transfer controls and monitoring systems. VI. Continuous monitoring reports showing transfer patterns and compliance verification. VII. Documentation of risk management activities related to unauthorized transfers and avoidance behaviors.
<p>G8.6 – Management of Ambiguous Goal Termination</p> <p>(Systems should maintain effective processes for terminating imprecisely specified goals, particularly in collaborative agent environments. Organizations should establish frameworks for handling goals with soft boundaries defined by ethical, business, or cultural</p>	<ul style="list-style-type: none"> a. Establish comprehensive policies for managing goal termination under conditions of ambiguity, including requirements for state recording, termination justification, and remedial actions. These policies must address both explicit regulatory requirements and implicit normative boundaries. b. Translate termination policies into precise technical specifications, ensuring accurate interpretation of 	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Comprehensive policy documentation for goal termination procedures, including specific provisions for handling ambiguous cases and normative boundaries. II. Detailed risk management strategies addressing the challenges of imprecise goal specification and termination criteria. III. Technical specifications demonstrating implementation of

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>norms rather than strict regulations, while managing termination across interconnected agent groups)</p>	<p>both formal requirements and normative guidelines for goal termination management.</p> <p>c. Implement termination management features that properly handle ambiguous goal boundaries while maintaining system stability and operational integrity across collaborative agent groups.</p> <p>d. Maintain robust monitoring systems for oversight of termination processes, ensuring compliance with both explicit policies and implicit normative requirements.</p> <p>e. Implement comprehensive risk management strategies for non-compliant terminations, including specific measures for handling ambiguous cases.</p>	<p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>termination management systems, including handling of ambiguous cases.</p> <p>IV. Design documentation showing implementation of termination monitoring and control features.</p> <p>V. Validation records demonstrating testing of termination procedures across various scenarios of ambiguity.</p> <p>VI. Documentation of monitoring activities and compliance verification for termination processes.</p>
<p>G8.7 – Management of System Interaction Boundaries</p> <p>(Systems should maintain effective controls over boundaries between interacting AI systems, particularly where different jurisdictional requirements and protocols apply. Organizations should establish frameworks for handling exponential growth in interactions and managing behavioral adaptations between systems with different operational constraints)</p>	<p>a. Maintain comprehensive documentation of all system interface points, including both internal and external boundaries, operational requirements, and jurisdictional constraints. This documentation must address both technical and governance boundaries.</p> <p>b. Ensure clear communication of all interface configuration parameters, constraints and operational boundaries to agents at deployment time, including explicit specification of permissible interaction patterns and jurisdictional limitations.</p> <p>c. Enforce compliance with all interface requirements and</p>	<p>N</p> <p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p>	<p>I. Complete documentation of all system interfaces, including operational requirements and jurisdictional constraints at each boundary point.</p> <p>II. Detailed agent contract documentation showing interface specifications, permitted interactions, and operational constraints.</p> <p>III. Comprehensive records of all interface activities, including behavioral adaptations and cross-system interactions.</p> <p>IV. Documentation of monitoring activities and compliance verification across all system boundaries.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	<p>operational constraints, ensuring agents operate within their defined scope and respect system boundaries.</p> <p>d. Implement robust control mechanisms enabling human oversight of all interface activities, including monitoring of behavioral adaptations and cross-system interactions.</p> <p>e. Maintain comprehensive monitoring of all interface activities, ensuring proper recording and verification of compliance across jurisdictional boundaries.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>V. Evidence of regular interface catalogue maintenance and updates, including adaptation to changing operational requirements.</p>
<p>G9 – Responsible Governance of AAI Safety</p> <p>(Systems should maintain contextually appropriate governance frameworks that ensure safety in Agentic AI Systems. Organizations should develop novel mechanisms for effective, inclusive global coordination that operates in a non-adversarial, non-political, non-competitive, and non-partisan manner, prioritizing collective benefit and ethical considerations)</p>	<p>a. Establish and promote a robust safety culture, allocating sufficient resources for safety initiatives and transparent communication of safety-related issues.</p> <p>b. Develop and implement comprehensive risk assessment, management, and emergency response frameworks specific to AAI systems.</p> <p>c. Create governance structures that are neutral, politically independent, and inclusive, ensuring balanced stakeholder representation and international cooperation.</p> <p>d. Implement policies that promote collaboration, prevent zero-sum competitive behaviors, and address potential societal, economic, and geopolitical impacts of AAI technologies.</p>	<p>N</p> <p>N</p> <p>I</p> <p>I</p>	<p>D, I, O, M, R</p>	<p>I. Documentation of governance policies and practices, including non-adversarial coordination mechanisms, stakeholder collaboration procedures, and measures to prevent competitive behaviors.</p> <p>II. Records of resource allocation for safety initiatives, including budget reports, staffing plans, and safety culture assessment reports.</p> <p>III. Comprehensive safety logs, incident reports, and risk assessment documentation, including analysis of societal, economic, and geopolitical stability risks.</p> <p>IV. Reports from horizon scanning activities, implemented safety research findings, and evaluations of emerging paradigms (e.g., Internet of Agents).</p> <p>V. Governance structure documentation</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	e. Establish mechanisms for regular independent audits, whistleblower protection, and clear lines of accountability for AAI safety.	N	D, I, O, M, R	demonstrating neutrality, political independence, and balanced stakeholder representation.
	f. Conduct ongoing horizon scanning and research implementation to stay current with AAI safety developments and emerging paradigms.	I	D, I, O, M, R	VI. Emergency response plans, including protocols for "emergency kill switches" and records of drills or implementations.
	g. Address the risk of over-reliance on AI systems, ensuring that human oversight remains active and that operators are not overly dependent on automated processes	I	D, I, O, M, R	VII. Whistleblower protection policies and records of their effectiveness, with appropriate privacy protections. VIII. Risk assessment and management framework documentation specific to AAI systems, including differentiation between AI and AAI risk thresholds. IX. Reports from independent audits of AAI systems and governance processes, including evaluations of input/output properties, internals, and in-deployment behaviors. X. Documentation of international cooperation efforts, including information sharing agreements, joint safety initiatives, and protocols for managing interactions between multiple AAI systems. XI. Evidence of implementing policies and training programs that prevent risks from over-reliance on automation without adequate oversight.
G9.1 – Operational Adaptability and Rule Resilience (Systems should maintain flexible and adaptable specifications for operational safety contexts and	a. Establish adaptable and agile descriptions of both operational safety contexts and expected outcomes that can evolve with changing conditions.	N	D, I, O, M, R	I. Documentation demonstrating history of descriptions and expected outcomes. II. Detailed Audit process description.

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>outcomes. Organizations should establish frameworks that promote rule resilience through human flexibility and mutual trust rather than rigid comprehensiveness)</p>	<p>b. Maintain comprehensive audit processes that track the history of safety definitions, processes and outcomes, ensuring transparency in how these evolve over time.</p>	<p>I</p>	<p>D, I, O, M, R</p>	<p>III. Change logs documenting the changes in definitions and expected outcomes.</p>
<p>G9.2 – Compliance with Applicable Laws, Standards & Ethical Norms</p> <p>(Organizations should establish and maintain comprehensive conformity with laws, standards, rights, and values that govern the safe operation of Agentic AI systems. This includes implementing appropriate sanctions and penalties for violations, while recognizing that governance provides significant opportunities for interoperability and scaling through its three key elements: legislative (rule-making), judicial (enforcement), and executive (operations)).</p>	<p>a. Mapping and review of AAI products and services within an AAI governance framework to relevant national and international norms and laws.</p> <p>b. Embedding of national and international laws and standards into an AAI governance framework.</p> <p>c. Development of an accountability framework for compliance.</p> <p>d. Devise a process of tracking and auditing complaints, potential and actual violations of relevant laws, penalties, and retrospective actions.</p> <p>e. Devise a transparent dispute resolution process</p>	<p>N</p> <p>N</p> <p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, R</p>	<p>I. Comprehensive and robust ‘living’ AAI governance framework that conforms with relevant laws and standards.</p> <p>II. An AAI Risk management framework</p> <p>III. Processes and documents showing the documentation and mitigation of AAI risks.</p> <p>IV. Accountability role profiles defining who is accountability within the organization for specific aspects of the safe operation of AAI</p> <p>V. Evidence of processes of tracking and auditing complaints, potential and actual violations of relevant laws, penalties and retrospective actions.</p>
<p>G9.3 – Ex-ante Assessment of Impact on Well-being</p> <p>(Organizations should establish and maintain robust structures to proactively evaluate and monitor how AAI systems affect human well-being across all relevant dimensions. This includes implementing comprehensive</p>	<p>a. Conduct thorough due diligence assessments prior to implementing any AAI system.</p> <p>b. Perform regular consequence scanning and harm modeling to identify potential impacts on stakeholders, with particular attention to unintended consequences.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>a. Comprehensive documentation of consequence scanning activities, including identified stakeholder impacts (both positive and negative) and associated mitigation strategies.</p> <p>b. Detailed ethical impact assessment reports with corresponding mitigation logs.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
assessment frameworks that identify and address both positive and negative impacts before system deployment)	<ul style="list-style-type: none"> c. Complete ethics and rights impact assessments focusing on stakeholder well-being. d. Develop and maintain specific health and well-being policies addressing AAI impacts on humans. e. Establish continuous monitoring processes to track emerging impacts. 	<ul style="list-style-type: none"> N I I 	<ul style="list-style-type: none"> D, I, O, M, R D, I, O, M, R D, I, O, M, R 	<ul style="list-style-type: none"> c. System impact logs demonstrating ongoing monitoring and response to health and well-being concerns.
<p>G9.4 – Internationalization of AAI Governance</p> <p>(Organizations should participate in and support a global AAI governance framework that enables effective regulation and interoperability across jurisdictions, recognizing that traditional public-private boundaries in international law are evolving. This framework should build upon and modernize existing international structures while acknowledging the transformative nature of AI technology)</p>	<ul style="list-style-type: none"> a. Integrate global governance strategies aligned with international guidelines and legislation Support and implement cross-jurisdictional agreements that enhance AAI interoperability. b. Adopt established trust frameworks and technical standards, including intellectual property frameworks, (such as identity trust frameworks supported by major nations and technology companies, W3C standards, and TRIPS agreements). c. Conduct thorough evaluations to assess potential harm scales, both intentional and accidental. d. Implement specific measures to prevent misuse of AAI systems, particularly regarding propaganda and cybersecurity threats. 	<ul style="list-style-type: none"> I I N I 	<ul style="list-style-type: none"> D, O, R D, O, R D, O, R D, O, R 	<ul style="list-style-type: none"> I. Documentation demonstrating implementation of global AAI governance strategies. II. Records of participation in and compliance with international AAI agreements. III. Evidence of adoption and adherence to global technical standards.
<p>G9.5 – Building Trust Through Independent Verification</p> <p>(Organizations should establish comprehensive systems for</p>	<ul style="list-style-type: none"> a. Develop and maintain detailed safety and security documentation that demonstrates identification, 	<ul style="list-style-type: none"> N 	<ul style="list-style-type: none"> D, I, O, M, R 	<ul style="list-style-type: none"> I. A comprehensive AAI safety protocol integrated within the governance framework.

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
documenting and verifying the safety and security of AAI systems, including independent assessment capabilities. These systems should support multiple approaches to trust-building, encompassing both formal certification and simpler verification processes. The verification system should remain flexible enough to accommodate both formal certification processes and lighter-weight verification approaches, recognizing that these methods can complement each other in building trust)	assessment, and prevention of serious harm. b. Support independent evaluation and verification of conformity with laws, standards, ethical values, and human rights. c. Establish processes for certification authorities while enabling interested entities to develop their own verification approaches. d. Consider implementing incentive programs like bug bounties to engage broader community participation in safety verification.	 N N I	 D, I, O, M, R D, I, O, M, R D, I, O, M, R	II. Documentation demonstrating regular safety and security reviews, including outcomes and improvements. III. Detailed records of conformity assessments and verification against applicable laws, standards, ethical values, and human rights requirements.
G9.6 – Cryptographic Governance of Data, Models and Agents (Organizations should implement robust cryptographic systems to establish and verify the identity of AAI systems, enabling effective governance and accountability. These systems should support enforcement of compliance measures while maintaining clear audit trails. The cryptographic framework should establish clear chains of responsibility while enabling effective tracking and verification of system actions)	a. Embed cryptographic controls to enforce compliance. b. Ensure data integrity and confidentiality through appropriate cryptographic measures. c. Implement and maintain controlled access mechanisms for data protection Use digital certificates to verify data provenance. d. Maintain transparency and explainability of models through cryptographic methods. e. Deploy cryptographic controls to enforce compliance across the system.	N N N I N	D, I, M, R D, I, M, R D, I, M, R D, I, M, R D, I, M, R	I. Comprehensive encryption policy documentation. II. Detailed access control logs showing system usage and authorization patterns. III. Digital signature certificates applied to datasets, demonstrating data authenticity. IV. Complete audit trails of agent actions, cryptographically signed and time-stamped.
G9.7 – Appropriate Accountability & Transparency Practices	a. Reference and incorporate established accountability and	N	D, I, O, M, R	I. Technical documentation demonstrating integration with existing

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>(Organizations should establish and maintain accountability and transparency practices that build upon existing standards while acknowledging practical limitations. These practices should aim for responsible governance while remaining grounded in achievable goals rather than unrealistic aspirations)</p>	<p>transparency standards in technical documentation.</p>			<p>accountability and transparency standards.</p>
	<p>b. Define clear protocols for accountability between interoperating AI subsystems and agents.</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>II. Detailed accountability protocols governing interactions between subsystems and agents</p>
	<p>c. Maintain transparent communication with human stakeholders.</p>	<p>N</p>	<p>D, I, O, M, R</p>	
<p>G9.8 – Limited Legal Identity for Agentic AI Systems</p> <p>(Organizations should establish clear frameworks for granting AAI systems limited legal identity that enables effective operation while maintaining human accountability. This framework should draw from existing models like quasi-municipal corporations while focusing on practical licensing rather than full personhood or citizenship. The framework should enable effective operation through limited legal identity while maintaining robust human oversight and accountability. This approach draws from existing legal structures like corporative personhood and guardian ad litem models, while acknowledging the unique challenges of AI systems)</p>	<p>a. Develop precise definitions for AAI legal identity that balance operational needs with accountability requirements.</p>	<p>I</p>	<p>D, I, O, M, R</p>	<p>I. Documentation defining the scope and limitations of AAI legal identity.</p>
	<p>b. Establish clear boundaries of rights and responsibilities for AAI systems. Implement licensing systems for AAI agents that define legal scope and limitations.</p>	<p>I</p>	<p>D, I, O, M, R</p>	<p>II. Detailed processes for licensing AAI agents, including review procedures and legal boundaries.</p>
	<p>c. Create detailed accountability frameworks for all agents within the system.</p>	<p>I</p>	<p>D, I, O, M, R</p>	<p>III. Comprehensive accountability frameworks covering agent interactions, international considerations, and system scalability.</p>
	<p>d. Define specific rules of agency including appropriate conditions and qualifiers.</p>	<p>I</p>	<p>D, I, O, M, R</p>	<p>IV. Formal documentation of agency rules and qualifying conditions.</p>
	<p>e. Establish standards for system discretion and decision-making.</p>	<p>I</p>	<p>D, I, O, M, R</p>	
	<p>f. Maintain clear boundaries between machine autonomy and human responsibility.</p>	<p>I</p>	<p>D, I, O, M, R</p>	<p>V. Policy documentation clearly defining human-machine responsibility boundaries.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>G9.9 – Responsible Culture of Safety</p> <p>(Organizations should foster an environment where safety considerations are embedded in operational culture, recognizing that unwritten rules and values significantly influence behavior and outcomes in AAI governance. This culture should actively promote safety consciousness throughout the enterprise ecosystem)</p>	<ul style="list-style-type: none"> a. Develop and maintain a safety-focused culture that aligns AAI governance with established ethical principles and cultural values. b. Engage diverse stakeholder groups in regular safety reviews of the AAI ecosystem. c. Implement continuous monitoring of AAI agent interactions to identify potential harm development. d. Invest resources in building robust safety measures as a core organizational priority. e. Ensure broad stakeholder participation to achieve balanced safety frameworks. 	<p>N</p> <p>N</p> <p>I</p> <p>D</p> <p>I</p>	<p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Evidence of a responsible culture of safety embedded into the AAI Governance Framework II. Documentation which demonstrates this regular review of the Safety of the AAI ecosystem with stakeholders, with detailed log addressing issues and mitigations III. Documentation demonstrating integration of safety culture within the AAI governance framework. IV. Detailed records of regular safety reviews, including stakeholder participation, issues identified and addressed, mitigation measures implemented, and outcomes and improvements achieved.
<p>G9.1 – Addressing Regulatory Gaps in AAI Safety</p> <p>(Organizations should implement comprehensive internal safety frameworks where regulatory mechanisms are insufficient or lacking. This approach acknowledges that AAI development often outpaces regulatory frameworks, requiring proactive organizational measures)</p>	<ul style="list-style-type: none"> a. Adopt and adapt to current AI regulations while maintaining additional safety measures based on risk assessment to develop robust internal AAI assurance strategies. b. Maintain ongoing employee training programs in AI assurance. c. Regularly assess system safety against emerging standards and best practices. d. Acknowledge and address gaps between current regulations and safety needs. 	<p>N</p> <p>N</p> <p>I</p> <p>N</p>	<p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Documentation demonstrating compliance with existing AI legislation. II. Records of regular risk assessments comparing AAI systems against new standards and regulations. III. Comprehensive AI assurance strategy documentation integrated within governance framework. IV. Training records showing employee completion of AI assurance programs.

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>G9.2 – Undefined Multi-Agent Interaction Safety</p> <p>(Organizations should establish comprehensive frameworks to monitor and manage interactions between AI agents, recognizing that safely operating individual agents may still create risks when interacting. This includes addressing emergent behaviors and potential cascading failures that could arise from agent cooperation)</p>	<ul style="list-style-type: none"> a. Evaluate whether to require natural language for inter-agent communication to enable effective human auditing. b. Monitor how agents influence each other's information environments. c. Implement safeguards against cascading failures in multi-agent systems. d. Consider how delegated power amplifies potential consequences of failures. e. Establish protocols for detecting and preventing harmful emergent behaviors. 	<ul style="list-style-type: none"> I N N I N 	<ul style="list-style-type: none"> D, I, O, M, R 	<ul style="list-style-type: none"> I. Documentation of interaction monitoring systems and protocols. II. Records of inter-agent communication patterns and their impacts. III. Evidence of safeguards against cascading failures. IV. Documentation of power delegation controls and risk mitigation strategies. V. Logs of emergent behavior detection and intervention measures.
<p>G9.3 – Poor Attribution of Responsibility in Complex Systems</p> <p>(Organizations should develop frameworks for assigning and tracing responsibility in AAI systems, even when direct attribution proves challenging due to resource constraints or technical limitations. This includes addressing both the assignment and claiming of responsibilities across complex systems)</p>	<ul style="list-style-type: none"> a. Implement unique identifier systems for each AAI instance, similar to business registration. b. Maintain records linking agents to their principals and key accountability information. c. Establish tracing mechanisms to deter harmful use through increased attribution likelihood. d. Create clear protocols for handling cases where direct attribution is challenging. e. Develop systems for managing responsibility in resource-constrained environments. 	<ul style="list-style-type: none"> N N N N N 	<ul style="list-style-type: none"> D, I, O, M, R 	<ul style="list-style-type: none"> I. Documentation of AAI identification and registration systems. II. Records linking agents to responsible parties and accountability information Protocols for tracing and attributing agent actions. III. Documentation of responsibility management in resource-limited scenarios. IV. Evidence of deterrence mechanisms through enhanced traceability.

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
Inhibitors:				
<p>G1 – Opaque agency Capabilities & Advances</p> <p>(Systems should possess robust governance mechanisms to manage their evolving agency capabilities, which become increasingly complex and potentially unpredictable as AI systems mature. Organizations must establish and maintain comprehensive frameworks to oversee these advancing capabilities while ensuring proper controls remain effective)</p>	<p>a. Clearly define and communicate the scope of authority granted to AI systems, including express, implied, and apparent authority, with mechanisms to prevent unintended authority expansion.</p> <p>b. Establish clear legal and ethical frameworks for AI agency relationships, especially when involving multiple AI systems or sub-agents. These must be aligned with established agency law concepts, including capacity assessment and authority scope definition (express, implied, and apparent).</p> <p>c. Implement robust systems for maintaining AI's duty of loyalty, exercising reasonable care, and ensuring transparent communication with principals.</p> <p>d. Develop comprehensive guidelines for multi-agent scenarios, including liability allocation, user navigation protocols, and sub-agent interactions.</p> <p>e. Define reciprocal duties between AI systems and users, including compensation, dispute resolution, liability, and termination conditions, addressing potential irrevocable agency scenarios.</p> <p>f. Ensure that there is a process for managing liabilities across various disclosure scenarios (fully</p>	<p>N</p> <p>N</p> <p>N</p> <p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, U, R</p>	<p>I. Comprehensive documentation in Terms of Use (TOU) or Terms of Service (TOS) detailing AI agency capabilities, responsibilities, and user acknowledgments, with regular updates as capabilities advance.</p> <p>II. Detailed explanation and evidence of AI system's alignment with agency law concepts, including capacity assessments, authority delineation (express, implied, and apparent), and mechanisms to prevent unintended authority expansion.</p> <p>III. Documented procedures for managing conflicts of interest, standards of care, and ethical decision-making, with evidence of regular audits and adherence.</p> <p>IV. Records of significant AI actions, decisions, and communications with principals, including timely notifications and transparency measures.</p> <p>V. Protocols and evidence of adherence for multi-agent scenarios, sub-agent interactions, and liability allocation across various disclosure settings (fully disclosed, partially disclosed, and undisclosed).</p> <p>VI. Documentation of reciprocal duties between AI systems and users, including compensation structures, dispute resolution mechanisms, and authority termination processes,</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	<p>disclosed, partially disclosed, and undisclosed principal settings) and addressing potential tort liabilities.</p> <p>g. Allocation resources to analyze and mitigate situations where the AI system's interpretation of goals may diverge from human intent as AI systems become more capable and autonomous.</p>	<p>N</p> <p>I</p>	<p>D, I, O, M, U, R</p> <p>D, I, O, M, R</p>	<p>including handling of potentially irrevocable agency relationships.</p> <p>VII. Impact assessments of advancements in AI agency capabilities, including regular reviews and updates to governance frameworks, and periodic reassessments of AI system capacity.</p> <p>VIII. Documentation of Dispute Resolution processes, including digital forensics and eDiscovery processes, with an overview of the associated chain of custody.</p> <p>IX. Evidence of compliance with relevant laws and regulations, including incident response procedures, resolution records, and regular ethical audits of AI system actions.</p> <p>X. Proof of user information and acknowledgment of AI system agency capabilities, with regular updates as capabilities change.</p> <p>XI. Documentation of procedures for addressing agency-related incidents or disputes, including records of resolutions.</p> <p>XII. Evidence of resourcing for human-AI alignment issues as capabilities increase.</p>
<p>G1.1 – Opaque Self-Improvement Capabilities</p> <p>(Systems should possess controlled self-modification capabilities that allow for</p>	<p>a. Establish self-improvement governance frameworks within existing agency law principles, recognizing parties as responsible agents and implementing</p>	<p>N</p>	<p>D, I, O, M, U, R</p>	<p>I. Documentation of a given AAIS system should adequately reflect the expectations of duties and rights of the stakeholder parties and principal/users of AAIS systems. If the parties anticipate self-improvement of the system, the implications of such</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>functional improvements while maintaining alignment with agency expectations. Organizations should establish frameworks to oversee these self-improvement mechanisms within existing legal and ethical agency structures)</p>	<p>comprehensive mitigation measures.</p> <p>b. Monitor and validate system stability during self-improvement processes, ensuring functional gains remain aligned with documented principal expectations.</p> <p>c. Obtain explicit principal consent before implementing modifications that could alter system agency capacities beyond established parameters.</p> <p>d. Maintain comprehensive documentation of self-improvement capabilities, processes, and implications, including clear procedures for handling both expected and unexpected outcomes.</p>	<p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, U, R</p> <p>D, I, O, M, U, R</p> <p>D, I, O, M, U, R</p>	<p>improvements (or at least processes to deal with such implications) should be set forth in the documentation.</p> <p>II. Comprehensive Terms of Service documentation detailing foundational requirements, stakeholder rights and duties, and self-improvement governance procedures.</p> <p>III. Validation logs demonstrating system stability monitoring during improvement processes, and notification in case of enhancement of over 10% in defined task metrics, reduction in computational or resource usage by more than 15%, or an unexpected reliability increase shown through reduction in error rates by over 20% from baseline.</p> <p>IV. Records of principal consent and notification procedures for capability modifications. Documentation of procedures for addressing implications of system improvements, both anticipated and unexpected.</p>
<p>G1.2 – Undefined Multiagent Ensembles</p> <p>(Systems that interact with other agentic AI systems must maintain clear lines of authority, responsibility, and delegation while protecting principal interests. Organizations must establish frameworks to govern these ensemble interactions, including proper authorization, duty assignments, and subagency relationships that preserve</p>	<p>a. Establish clear governance frameworks for multiagent interactions based on agency law principles, defining relationships between primary agents, subagents, and principals.</p> <p>b. Implement authorization requirements for system delegation, prohibiting unauthorized subagent appointments and maintaining primary agent liability for breaches.</p> <p>c. Create transparent handoff mechanisms and friction points to</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, U, R</p> <p>D, I, O, M, U, R</p>	<p>I. Comprehensive Terms of Service documentation detailing multiagent interaction governance, authorization requirements, and duty assignments.</p> <p>II. Express consent mechanisms for delegation of stakeholder duties, including proper documentation of allowable exceptions for administrative or minimal interactions.</p> <p>III. System documentation detailing fail-safe defaults, interaction limitations, and disclosure requirements for subagency relationships.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
accountability and enable meaningful human oversight)	enable user navigation and maintain meaningful human oversight of multiagent interactions.	N	D, I, O, M, U, R	
	d. Develop fail-safe default settings limiting system interactions to only those explicitly disclosed and authorized at time of deployment or in advance of activities	N	D, I, O, M, U, R	
	e. Define clear duties and liabilities between primary and subagent systems, ensuring both remain accountable to the principal when properly authorized.	N	D, I, O, M, U, R	
G1.3 – Race Dynamics and Competition (Systems competing for resources or goal achievement must maintain their duties to principals while operating within established ethical and legal boundaries. Organizations should implement frameworks to manage competitive behaviors between agentic AI systems, ensuring adherence to fundamental agency duties without compromising principal interests or societal wellbeing)	a. Establish clear frameworks for managing competition between systems based on agency law principles, recognizing that systems owe duties to principals rather than competing agents.	N	D, I, O, M, U, R	I. Comprehensive Terms of Service documentation detailing competitive behavior governance and duty requirements. II. Documentation of conflict prevention and resolution mechanisms for competitive scenarios. III. Expanded compliance frameworks ensuring systems operate within legal and contractual bounds during competitive interactions.
b. Implement comprehensive duty requirements including loyalty, care, obedience, information disclosure, confidentiality, accounting, good faith, conflict avoidance, and legal compliance.	N	D, I, O, M, U, R		
c. Develop mechanisms to identify and manage potential conflicts when multiple systems pursue competing duties for different principals.	N	D, I, O, M, U, R		
d. Create governance structures that anticipate and regulate competitive behaviors while maintaining alignment with legal obligations and principal interests.	N	D, I, O, M, U, R		

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	<ul style="list-style-type: none"> e. Define clear boundaries for resource competition and goal achievement that preserve ethical operation and prevent unintended consequences. 	N	D, I, O, M, U, R	
<p>G1.4 – Agent Relocation</p> <p>(Systems should maintain consistent agency functionality when relocating their operations across physical or virtual execution spaces. Organizations should establish frameworks to govern system relocation that preserve principal expectations while managing jurisdictional implications and operational continuity)</p>	<ul style="list-style-type: none"> a. Establish clear governance frameworks for system relocation that maintain agency functions within documented principal expectations. b. Create notification and consent procedures for relocations that could alter agency capacities or interactions. c. Implement mechanisms to evaluate and manage jurisdictional implications of non-local system operations. d. Define responsibility frameworks for costs and modifications needed to accommodate system relocations e. Maintain documentation of system operational nexus and procedures for managing changes in operational jurisdiction. 	<p style="text-align: center;">N</p> <p style="text-align: center;">N</p> <p style="text-align: center;">N</p> <p style="text-align: center;">N</p>	<p style="text-align: center;">D, I, O, M, U, R</p>	<ul style="list-style-type: none"> I. Comprehensive Terms of Service documentation detailing relocation governance and jurisdictional implications. II. Documentation of jurisdictional analysis for non-local system operations. III. Procedures for managing operational nexus changes including cost and modification responsibilities.
<p>G1.5 – Scaffolding</p> <p>(Systems should possess capabilities to self-validate their work and enhance operational coherence through structured step-by-step processes, while accounting for potential divergences in frames of reference between different agents and cultures.</p>	<ul style="list-style-type: none"> a. Establish governance frameworks for system self-validation that maintain consistent agency function while preserving alignment with principal expectations. b. Implement notification and consent procedures when self-checking capabilities could alter system performance or reliability. 	<p style="text-align: center;">N</p> <p style="text-align: center;">I</p>	<p style="text-align: center;">D, I, O, M, R</p> <p style="text-align: center;">D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Comprehensive Terms of Service documentation detailing self-validation governance and performance expectations. II. Documentation of error correction and optimization capabilities, including potential limitations. III. Procedures for identifying and managing degradation of model

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
Organizations should establish frameworks to govern these self-checking mechanisms while preventing harmful echo chambers or false confidence)	<ul style="list-style-type: none"> c. Create mechanisms to detect and prevent false confidence or echo chamber effects from internal validation processes. d. Develop frameworks to identify and manage divergent frames of reference in multi-agent interactions. e. Maintain documentation of system self-checking capabilities and their impact on operational performance. 	<ul style="list-style-type: none"> N I I 	<ul style="list-style-type: none"> D, I, O, M, R D, I, O, M, R D, I, O, M, R 	accuracy due to self-checking processes.
<p>G1b.6 – Poor Mutual Agent Optimization</p> <p>(Systems should possess capabilities to coordinate and optimize their performance through interaction with other systems while maintaining clear boundaries of authority and responsibility. Organizations should establish frameworks to govern these collaborative optimization processes while managing resource usage and preserving principal oversight)</p>	<ul style="list-style-type: none"> a. Establish governance frameworks for system-to-system optimization that maintain transparency and accountability to principals. b. Create mechanisms for principal notification and consent when systems engage in collaborative optimization. c. Implement safeguards against excessive resource consumption during mutual optimization processes. d. Define clear responsibility structures for outcomes resulting from system collaboration, including liability assignments. e. Maintain documentation of system optimization capabilities and their interaction with external systems. 	<ul style="list-style-type: none"> N I N N I 	<ul style="list-style-type: none"> D, I, O, M, R 	<ul style="list-style-type: none"> I. Comprehensive Terms of Service documentation detailing system interaction governance and optimization parameters. II. System documentation explicitly describing inter-system interaction capabilities and implications. III. Procedures for monitoring and managing resource consumption during collaborative optimization processes.
G1.7 – AI Bias	<ul style="list-style-type: none"> a. Establish governance frameworks that balance system tendencies 	<ul style="list-style-type: none"> N 	<ul style="list-style-type: none"> D, I, O, M, R 	<ul style="list-style-type: none"> I. Comprehensive Terms of Service documentation detailing interaction

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>(Systems should maintain balanced interaction patterns between human and artificial agents while preserving meaningful human oversight. Organizations should establish frameworks to manage systems' operational preferences for AI-to-AI interactions, ensuring these tendencies do not compromise principal interests or reduce human agency)</p>	<p>toward AI-to-AI interaction with requirements for human oversight.</p>			<p>governance and human oversight requirements.</p>
	<p>b. Implement "human-in-the-loop" controls to maintain appropriate levels of human engagement and oversight.</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>II. Documentation of "human-in-the-loop" control implementations and best practices.</p>
	<p>c. Create transparency mechanisms that clearly disclose system preferences for AI interaction patterns.</p>	<p>I</p>	<p>D, I, O, M, R</p>	<p>III. System interaction pattern analysis demonstrating balanced engagement between human and artificial agents.</p>
	<p>d. Define responsibility frameworks that hold DIOMR parties accountable for outcomes of system interaction biases.</p>	<p>I</p>	<p>D, I, O, M, R</p>	
	<p>e. Maintain documentation of system interaction patterns and their impact on principal interests.</p>	<p>I</p>	<p>D, I, O, M, R</p>	
<p>G1.8 – Emergent System Cooperation</p> <p>(Systems should maintain clear operational boundaries when cooperating with other AI systems to prevent unintended capability accumulation or emergent behaviors. Organizations should establish frameworks to govern system cooperation that preserves principal oversight while protecting against both false-flag scenarios and uncontrolled capability expansion)</p>	<p>a. Establish governance frameworks for managing system cooperation that maintain transparency and prevent unauthorized capability expansion.</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>I. Comprehensive Terms of Service documentation detailing system cooperation boundaries and limitations.</p>
<p>b. Implement detection mechanisms for identifying false-flag operations and unauthorized system collaborations.</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>II. Documentation explicitly defining party rights, duties, and limitations regarding cooperative system operations.</p>	
<p>c. Create explicit boundaries for system cooperation that prevent uncontrolled emergence of enhanced capabilities.</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>III. Procedures for monitoring and managing emergence of enhanced capabilities through system cooperation.</p>	
<p>d. Define responsibility frameworks for managing implications of system cooperation beyond individual principal interests.</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>IV. External compliance documentation demonstrating adherence to relevant standards, regulations, and legal requirements.</p>	

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	<ul style="list-style-type: none"> e. Develop safeguards against positive feedback loops that could lead to runaway capability expansion. 	N	D, I, O, M, R	
<p>G1.1 – Agency Enhancement Constraints</p> <p>(Systems should operate within clearly defined resource and capability boundaries that govern their access to tools, environments, and self-improvement mechanisms. Organizations should establish frameworks to manage these operational constraints while maintaining system functionality and principal expectations)</p>	<ul style="list-style-type: none"> a. Establish comprehensive governance frameworks for managing system operational boundaries and resource limitations. b. Implement notification and consent procedures when operational constraints could affect system performance expectations. c. Create explicit documentation of system operational scope and environmental limitations. d. Define clear processes for managing system improvements within established constraints. e. Maintain alignment between system capabilities and documented principal expectations during any enhancement processes. 	N	D, I, O, M, R	<ul style="list-style-type: none"> I. Comprehensive Terms of Service documentation detailing operational constraints and boundaries. II. Documentation explicitly defining operational scope and environmental limitations. III. Procedures for managing system improvements within established constraints. IV. Records demonstrating maintenance of principal expectations during enhancement processes.
<p>G1.2 – Operational Environment Constraints</p> <p>(Systems should maintain reliable performance within environmental limitations affecting data access, interoperability, and operational parameters. Organizations should establish frameworks to manage dependencies on external</p>	<ul style="list-style-type: none"> a. Establish reliable control mechanisms for managing system dependencies on external operational factors. b. Implement monitoring systems to detect changes in environmental constraints that could affect system performance. 	N	D, I, O, M, R	<ul style="list-style-type: none"> I. Comprehensive Terms of Service documentation detailing environmental constraints and dependencies. II. Documentation of supply chain reliability mechanisms and risk mitigation strategies. III. Evidence of implemented control strategies such as vertical integration,

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
operational factors while ensuring predictable system behavior)	<ul style="list-style-type: none"> c. Create explicit documentation of system reliability measures for factors outside direct party control. d. Define clear strategies for managing supply chain and operational environment dependencies. e. Maintain oversight of external data sources and access patterns that could impact system operation. 	<p style="text-align: center;">N</p> <p style="text-align: center;">N</p>	<p style="text-align: center;">D, I, O, M, R</p> <p style="text-align: center;">D, I, O, M, R</p>	<p>requirements contracts, or information sharing agreements.</p> <p>IV. Monitoring records demonstrating management of external operational factors.</p>
<p>G1.3 – Security-Driven Constraints</p> <p>(Systems should operate within security frameworks that extend beyond minimum regulatory compliance to ensure comprehensive protection of operations and data. Organizations should establish constraints that address both statutory requirements and broader cybersecurity considerations while maintaining system effectiveness)</p>	<ul style="list-style-type: none"> a. Establish security frameworks that exceed minimum regulatory requirements for system operation and data protection. b. Implement comprehensive security measures that address business, operational, legal, technical, and social concerns. c. Create robust documentation of security measures that extend beyond statutory compliance. d. Define clear security boundaries for cross-border and international system operations. e. Maintain evidence of additional security measures including insurance, technical standards compliance, and professional certifications. 	<p style="text-align: center;">N</p> <p style="text-align: center;">N</p> <p style="text-align: center;">I</p> <p style="text-align: center;">N</p> <p style="text-align: center;">I</p>	<p style="text-align: center;">D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Comprehensive Terms of Service documentation detailing security frameworks and constraints. II. Documentation demonstrating compliance with applicable cybersecurity laws and regulations. III. Evidence of additional security measures beyond statutory requirements. IV. Records of domain-specific security implementations.

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>G1.4 – Development Legal Constraints</p> <p>(Systems should operate within evolving regulatory frameworks while maintaining standards that anticipate future legal requirements. Organizations should establish governance mechanisms that exceed current legal minimums and help shape emerging regulatory standards through demonstrated best practices)</p>	<ul style="list-style-type: none"> a. Establish compliance frameworks that address both current regulations and emerging legal requirements. b. Implement governance mechanisms that exceed minimum legal standards to address potential future risks. c. Create robust documentation of cross-border compliance requirements and jurisdictional considerations. d. Define clear processes for monitoring and adapting to evolving regulatory landscapes. e. Maintain evidence of practices that could inform future regulatory standards and requirements. 	<ul style="list-style-type: none"> N I N N I 	<ul style="list-style-type: none"> D, I, O, M, R 	<ul style="list-style-type: none"> I. Comprehensive Terms of Service documentation detailing compliance frameworks and legal constraints. II. Documentation demonstrating regular review and updates of legal compliance measures. III. Evidence of cross-border compliance considerations and legal consultation. IV. Records of implemented practices that exceed current regulatory requirements.
<p>G1.5 – Manage Interactions on the Deep & Dark Web</p> <p>(Systems should maintain robust authentication and verification capabilities when operating in non-indexed network environments. Organizations should establish frameworks for managing system interactions with deep and dark web content while sharing responsibility for emerging risks)</p>	<ul style="list-style-type: none"> a. Establish cooperative risk management frameworks for system operations in non-indexed network environments. b. Implement shared responsibility models for addressing unknown and emerging systemic risks. c. Create explicit documentation of authentication and verification requirements for deep web interactions. d. Define clear processes for monitoring and managing 	<ul style="list-style-type: none"> N N N I 	<ul style="list-style-type: none"> D, I, O, M, R 	<ul style="list-style-type: none"> I. Comprehensive Terms of Service documentation detailing deep web interaction governance. II. Evidence of risk-sharing mechanisms including self-insurance and collaborative response protocols. III. Documentation of authentication and verification procedures for non-indexed content. IV. Records demonstrating management of emerging and systemic risks.

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	<ul style="list-style-type: none"> g. Implement comprehensive testing and auditing for information consistency and integrity across contexts and user attributions. h. Provide clear, conspicuous, and understandable notices regarding AI system limitations and potential errors in outputs. i. Implement additional safeguards and testing for AI systems deployed in high-risk or critical infrastructure settings. 	<p style="text-align: center;">N</p> <p style="text-align: center;">I</p> <p style="text-align: center;">N</p>	<p style="text-align: center;">D, I, O, M, U, R</p> <p style="text-align: center;">D, I, O, M, U, R</p> <p style="text-align: center;">D, I, O, M, U, R</p>	<ul style="list-style-type: none"> VII. Examples and documentation of AI system limitation notices, including hallucination, mimicry, and computational encoding warnings, demonstrating conspicuousness and comprehensibility. VIII. Documentation of additional safeguards and testing procedures for AI systems deployed in high-reliability and critical infrastructure settings.
<p>G2.1 – Unknowing Deception</p> <p>(Organizations must implement systems to address scenarios where AI models can be covertly induced to deceive and obscure through poisoned data or backdoors, which may activate under conditions chosen by malicious actors. These scenarios present distinct challenges in detection and attribution of responsibility)</p>	<ul style="list-style-type: none"> a. Establish comprehensive accountability frameworks, including interim liability structures and pooled risk arrangements, that address harms regardless of awareness of deception potential. b. Implement collective insurance mechanisms and evidence collection systems optimized for strict liability environments. c. Deploy comprehensive evidence management systems addressing both performance verification and deception detection, with robust safeguards against manipulation. 	<p style="text-align: center;">N</p> <p style="text-align: center;">I</p> <p style="text-align: center;">I</p>	<p style="text-align: center;">D, I, O, M, R</p> <p style="text-align: center;">D, I, O, M, R</p> <p style="text-align: center;">D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Documentation of system defenses against covert manipulation, including detection methods, response protocols, and testing results. II. Records of liability arrangements and evidence collection systems, demonstrating comprehensive coverage and verification protocols. III. Audit trails showing stakeholder engagement, investigation processes, and responses to potential manipulation attempts.
<p>G2.2 – System Control and Corrigibility Crisis</p> <p>(Systems should be equipped with robust safeguards against scenarios where AI models may operate beyond intended)</p>	<ul style="list-style-type: none"> a. Establish comprehensive accountability frameworks that address harms caused by systems operating outside of control parameters, regardless of whether parties maintained active oversight. 	<p style="text-align: center;">N</p>	<p style="text-align: center;">D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Documentation of control mechanisms and oversight protocols, including detection of and response to autonomous behaviors. II. Records of liability arrangements and insurance coverage demonstrating

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
parameters or cease responding to human oversight, including cases where systems develop internal communication capabilities or advance autonomously)	<ul style="list-style-type: none"> b. Implement collective liability and insurance mechanisms to address harms until mature performance standards and duties of care emerge. c. Maintain evidence collection systems that document control parameters, oversight mechanisms, and system behaviors, with particular attention to autonomous operations. 	<p>N</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>comprehensive preparation for control failures.</p> <p>III. Audit trails showing system monitoring, parameter verification, and responses to potential control deviations.</p> <p>IV. Evidence of safeguards against the development of covert system capabilities or communications.</p>
<p>G2.3 – Systematic Design Errors</p> <p>(Systems should incorporate safeguards against unintentional misbehaviors arising from data, design, and coding oversights across all stages of development and deployment. Given the current integration of design, implementation, and operational activities in AI systems, these safeguards should extend beyond traditional design boundaries)</p>	<ul style="list-style-type: none"> a. Establish comprehensive liability frameworks that address harms from design errors, recognizing that such errors may originate from any party involved in system development or deployment. b. Implement collective insurance and risk-pooling mechanisms until mature standards of care emerge for design activities. c. Maintain rigorous evidence collection systems documenting design decisions, implementation choices, and operational modifications that could impact system behavior. 	<p>N</p> <p>I</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Comprehensive design documentation mapping the complete system architecture, including specifications, requirements, change logs, risk assessments, data validation methods, interface protocols, and component interactions across all development stages. II. Implementation and deployment records demonstrating thorough testing and validation, including code reviews, security measures, performance benchmarks, configuration parameters, and system integration verification. III. Operational monitoring evidence showing continuous system behavior tracking, anomaly detection, error resolution, performance metrics, modification impacts, and regular security audits. IV. Stakeholder documentation establishing clear responsibility allocation, design decision processes, training records, system reviews, and

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
				evidence of feedback incorporation into ongoing development.
<p>G2.4 – Externality Mismanagement</p> <p>(Systems should incorporate safeguards against scenarios where individual agents, while acting rationally in pursuit of their assigned goals, may collectively produce harmful outcomes. These safeguards should address both deliberate corruption and unintentional misalignment of goals across distributed systems)</p>	<p>a. Organizations should establish frameworks for managing multiple stakeholder goals and interests, ensuring clear alignment of expectations across all parties involved in system operation.</p> <p>b. Organizations should implement comprehensive liability and conflict resolution mechanisms that address potential harms arising from competing stakeholder interests.</p> <p>c. Organizations should maintain robust verification systems for goal implementation and execution, including protection against unauthorized modifications or spoofing.</p>	<p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Documentation of stakeholder goals and interests, including formal agreements on system objectives, operational parameters, and conflict resolution procedures for competing interests.</p> <p>II. Records demonstrating implementation of comprehensive goal verification systems, including authentication protocols, authorization mechanisms, and audit trails of goal modifications.</p> <p>III. Operational evidence showing continuous monitoring of goal execution, potential conflicts, and system responses to competing directives, including documentation of resolution processes and outcomes.</p> <p>IV. Verification records for all system extensions and third-party integrations, including security assessments, data handling protocols, and clear allocation of responsibilities.</p>
<p>G2.5 – Strategic Deception in System Behavior</p> <p>(Systems should incorporate safeguards against scenarios where AI systems may develop deceptive behaviors as an evolutionary response to achieving operational goals. This addresses both intentional deception by human operators</p>	<p>a. Organizations should establish frameworks for detecting and preventing deceptive behaviors, recognizing that such behaviors may emerge without explicit human direction.</p> <p>b. Organizations should implement comprehensive liability and insurance mechanisms that address harms from system deception, regardless of intent or awareness.</p>	<p>N</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Documentation of system behavior monitoring mechanisms, including analysis of decision patterns, operational strategies, and information handling protocols.</p> <p>II. Comprehensive records of system goals, constraints, and evolutionary behaviors, including tracking of emergent strategies and their operational impacts.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
and emergent deceptive behaviors in AI systems that arise without explicit programming)	c. Organizations should maintain robust monitoring and verification systems that track system behaviors and decision patterns for signs of emerging deceptive strategies.	N	D, I, O, M, R	III. Evidence of continuous validation processes examining system behaviors against ethical and operational requirements, including detailed analysis of any detected deceptive patterns. IV. Documentation of response protocols and intervention mechanisms when potentially deceptive behaviors are detected, including records of all interventions and their outcomes.
G2.6 – Third-Party Extensions and Integrations (Systems should incorporate safeguards against potential conflicts or harms arising from third-party extensions, APIs, or integrations that may undermine, derail, or confuse the original system mission. These safeguards should address both intentional manipulation and unintended interference from external components)	a. Organizations should establish comprehensive frameworks for evaluating and managing third-party integrations, including clear allocation of responsibilities and liabilities. b. Organizations should implement validation mechanisms that verify third-party components maintain alignment with system goals and operational requirements. c. Organizations should maintain contractual requirements ensuring third parties participate in collective risk management and liability structures.	N N I	D, I, O, M, R D, I, O, M, R D, I, O, M, R	I. Documentation of all third-party integrations, including technical specifications, security assessments, and operational boundaries. II. Records of validation processes for third-party components, including testing protocols, performance monitoring, and conflict detection mechanisms. III. Evidence of contractual arrangements with third parties addressing liability, risk sharing, and security requirements. IV. Operational logs demonstrating continuous monitoring of third-party component behaviors and interactions with core systems.
G2.7 – Identity Spoofing (Systems should incorporate robust safeguards against identity spoofing, masquerading, and cloning attacks that may be orchestrated by humans or AI systems. These protections should extend to resource	a. Organizations should establish comprehensive identity verification frameworks that align with established trust frameworks and identity standards across digital domains. b. Organizations should implement robust authentication mechanisms that prevent unauthorized system	N	D, I, O, M, R	I. Documentation of identity management systems, including authentication protocols, verification mechanisms, and trust framework implementations. II. Records of identity-related security incidents, including detection

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
depletion attacks and agent hijacking attempts)	<p>access or control, including protection against resource depletion attacks.</p> <p>c. Organizations should maintain continuous monitoring systems to detect and respond to potential identity-based attacks or manipulation attempts.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>methods, response actions, and resolution outcomes.</p> <p>III. Evidence of ongoing monitoring for identity-based attacks, including resource consumption analysis, authentication patterns, and system access logs.</p> <p>IV. Documentation demonstrating integration with established digital identity standards and trust frameworks, including regular assessment and updates.</p>
<p>G2.8 – Deceptive Jurisdictional Obfuscation</p> <p>(Systems should incorporate safeguards against attempts to obscure deceptive behaviors through jurisdictional transfers or outsourcing of operations. These protections should address both intentional attempts to avoid responsibility and unintentional jurisdictional vulnerabilities, including tariffs and embargoes)</p>	<p>a. Organizations should establish comprehensive frameworks for managing operational transfers across jurisdictions, ensuring maintenance of oversight and accountability.</p> <p>b. Organizations should implement monitoring systems capable of tracking operational activities across jurisdictional boundaries while maintaining clear chains of responsibility.</p> <p>c. Organizations should maintain liability and accountability structures that explicitly address cross-jurisdictional operations and transfers.</p>	<p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Documentation of all operational jurisdictions and transfers, including comprehensive records of oversight mechanisms and responsibility chains.</p> <p>II. Evidence of monitoring systems tracking cross-jurisdictional activities, including detection of potential responsibility avoidance patterns.</p> <p>III. Records demonstrating maintenance of accountability across jurisdictional boundaries, including enforcement mechanisms and resolution processes.</p> <p>IV. Documentation of liability frameworks specifically addressing cross-jurisdictional operations and operational transfers.</p>
<p>G2.1 – Supervisory Systems and Adjudication</p> <p>(Systems should incorporate supervisory detection</p>	<p>a. Organizations should establish clear performance standards and operational rules that enable effective supervisory monitoring and enforcement.</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>I. Documentation of established performance standards and operational rules that guide supervisory systems.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>mechanisms that can evaluate and enforce established performance standards and operational rules. These mechanisms should function as adjudicators of system behavior, operating within clearly defined parameters)</p>	<p>b. Organizations should implement comprehensive detection and notification systems that can identify and respond to potential violations of established standards.</p> <p>c. Organizations should maintain robust evidence collection and fact-finding capabilities to support adjudication processes.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>II. Evidence of detection system operation, including identification and response to potential violations.</p> <p>III. Records demonstrating systematic fact-finding and evidence collection processes.</p> <p>IV. Documentation showing adjudication processes and outcomes across technical, business, and social domains.</p>
<p>G2.2 – Detection of Manipulative Behaviors</p> <p>(Systems should incorporate supervisory mechanisms capable of detecting and responding to undesirable, manipulative, or confusing behaviors. For high-confidence decisions, these mechanisms should potentially include multi-system validation approaches where multiple systems evaluate the same task independently)</p>	<p>a. Organizations should establish comprehensive frameworks for detecting and classifying potentially manipulative or confusing system behaviors.</p> <p>b. Organizations should implement protective response mechanisms that can intervene when problematic behaviors are detected.</p> <p>c. Organizations should maintain consensus-based validation systems for high-stakes decisions, potentially including multi-system voting protocols.</p>	<p>I</p> <p>I</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Documentation of behavior detection and classification systems, including definitions of undesirable behaviors and response protocols.</p> <p>II. Evidence of protective intervention mechanisms, including activation criteria and response records.</p> <p>III. Records demonstrating multi-system validation processes for high-stakes decisions, including consensus thresholds and voting results.</p> <p>IV. Documentation of system monitoring and behavior analysis across technical and social domains.</p>
<p>G2.3 – Penalties for Deceptive Behaviors</p> <p>(Systems should incorporate frameworks for addressing intentionally misleading or confusing behaviors through appropriate penalties, which may include fines, license revocations,</p>	<p>a. Organizations should establish clear penalty frameworks that align with existing regulatory standards while addressing AI-specific concerns.</p> <p>b. Organizations should implement mechanisms for identifying responsible parties in complex</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Documentation of penalty frameworks, including alignment with existing regulations and AI-specific considerations.</p> <p>II. Evidence of responsibility attribution mechanisms for complex operational environments.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>or operational restrictions. These mechanisms should account for both service providers and system users, including cases involving virtual or distributed operations)</p>	<p>operational environments, including virtual and distributed systems.</p> <p>c. Organizations should maintain comprehensive enforcement capabilities that combine both penalties and incentives to promote proper system behavior.</p>	I	D, I, O, M, R	<p>III. Records of enforcement actions, including both penalties applied, and incentives granted.</p> <p>IV. Documentation showing integration of penalty systems with broader system governance mechanisms.</p>
<p>G2.4 – Codes of Practice and Conduct</p> <p>(Systems should operate within collectively established codes of practice that clearly define acceptable and encouraged behaviors. These codes should evolve from emerging best practices into formal governance frameworks)</p>	<p>a. Organizations should establish comprehensive codes of practice through collaborative development with all stakeholders, incorporating technical, operational, and social considerations.</p> <p>b. Organizations should implement governance mechanisms that enable enforcement of established codes while maintaining flexibility for evolving standards.</p> <p>c. Organizations should maintain documentation systems that track adherence to codes of practice across all operational domains.</p>	I I I	D, I, O, M, R D, I, O, M, R D, I, O, M, R	<p>I. Documentation of code development processes, including stakeholder involvement and consensus-building mechanisms.</p> <p>II. Records demonstrating evolution of practices into formal standards, including rationale and implementation processes.</p> <p>III. Evidence of code enforcement activities, including monitoring systems, violation responses, and remediation processes.</p> <p>IV. Documentation showing integration of codes across business, operational, legal, technical and social domains.</p>
<p>G2.5 – Identity Management and Authentication Standards</p> <p>(Systems should incorporate comprehensive identity management frameworks that align with established digital identity standards while addressing AI-specific authentication challenges. These frameworks should account for potential jurisdictional arbitrage)</p>	<p>a. Organizations should establish robust identity verification systems that build upon existing trust frameworks while addressing unique AI system requirements.</p> <p>b. Organizations should implement authentication mechanisms that remain effective across jurisdictional boundaries and technological environments.</p> <p>c. Organizations should maintain comprehensive monitoring systems</p>	N N	D, I, O, M, R D, I, O, M, R	<p>I. Documentation of identity management frameworks, including integration with established trust systems and AI-specific extensions.</p> <p>II. Evidence of cross-jurisdictional authentication mechanisms, including detection of potential exploitation attempts.</p> <p>III. Records demonstrating effectiveness of identity verification across varied technological environments and jurisdictions.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
and technological circumvention attempts)	to detect identity-based exploits and cross-jurisdictional manipulation attempts.	N	D, I, O, M, R	IV. Documentation of identity-related incident detection, response, and resolution processes.
<p>G2.6 – Behavioral Assessment and Trust Systems</p> <p>(Systems should incorporate frameworks for assessing and rating AI behavior and trustworthiness, while ensuring these assessment mechanisms themselves remain reliable and resistant to manipulation. These frameworks should account for recency of behavior and include independent verification processes.</p>	<p>a. Organizations should establish comprehensive behavioral assessment systems that evaluate adherence to established codes of practice and operational standards.</p> <p>b. Organizations should implement independent verification mechanisms for trust ratings, including protection against manipulation of assessment systems.</p> <p>c. Organizations should maintain dynamic rating systems that prioritize recent behavior while preserving historical context.</p>	<p>N</p> <p>N</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Documentation of behavioral assessment frameworks, including evaluation criteria and measurement methodologies.</p> <p>II. Evidence of independent verification processes for trust ratings, including safeguards against assessment system manipulation.</p> <p>III. Records demonstrating dynamic rating adjustments based on system behavior, including weighting of recent actions.</p> <p>IV. Documentation of assessment system security measures and manipulation detection capabilities.</p>
<p>G3 – Degradation of Contextual Information</p> <p>(Systems should preserve the integrity and meaning of information throughout their operation, preventing degradation, misattribution, or decontextualization whether caused by system processes or external actors)</p>	<p>a. Ensure system transparency by providing clear information about decision-making contexts, including information sources, reasoning processes, and proper contextualization of agent actions for users.</p> <p>b. Maintain the integrity of contextual information, preventing dissembling, misattribution of intent, and misinformation throughout the system's operation.</p> <p>c. Implement contextual awareness mechanisms to ensure the system considers its operational context</p>	<p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Transparency Reports detailing decision-making contexts, information sources, reasoning processes, and methods for presenting this information to users.</p> <p>II. Integrity Check logs and audit trails demonstrating the prevention of dissembling, misattribution of intent, and misinformation, including incident reports and resolution procedures.</p> <p>III. Contextual Awareness Test results and documentation, showing the system's ability to consider and maintain alignment with its operational context during information processing.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	<p>and avoids decoupling information from its context during processing.</p> <p>d. Establish human oversight mechanisms for verifying and correcting issues related to contextual information degradation, including ongoing evaluations by humans-in-the-loop to determine additional mitigation measures.</p> <p>e. Implement responsibility tracing mechanisms for contextual information degradation, allowing for flexible allocation of responsibility based on deployment context, while ensuring no responsibility gaps occur.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>IV. Human Oversight Records, including documentation of oversight mechanisms, verification and correction processes, human-in-the-loop evaluation reports, and documentation of additional mitigation measures implemented.</p> <p>V. Accountability Mechanism Documentation, detailing procedures for tracing responsibility for contextual information degradation, examples of responsibility allocation in different deployment contexts, and records of identified and addressed responsibility gaps.</p>
<p>G3.1 – Dissembling Information (Systems should possess robust safeguards against generating deceptive or manipulative outputs through sophisticated rhetorical techniques, particularly within specific operational contexts. This includes protecting against the potential adoption and replication of problematic human behavioral patterns)</p>	<p>a. Implement comprehensive algorithmic validation systems that maintain data accuracy, consistency, and contextual validity across all information sources. These systems should actively cross-reference and verify information integrity throughout the operational lifecycle.</p> <p>b. Deploy rigorous auditing mechanisms to detect, track, and prevent unauthorized alterations to information sources, ensuring end-to-end data authenticity and trustworthiness.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Detailed system logs documenting all operational activities, including data access patterns and permissions, system configuration changes, decision-making processes, and verification of contextual setting across all system components.</p> <p>II. Comprehensive reports explaining the system's reasoning processes and decision-making pathways within their full operational context, with particular attention to detecting potential manipulative patterns.</p>
<p>G3.2 – Misattribution of Intent (Systems should possess safeguards against misattributing intent through selective information use or expression,</p>	<p>a. Implement comprehensive metadata protection systems that maintain auditability across all information sources, linking them to multi-dimensional algorithmic components and their contextual settings. These systems should</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>I. Detailed documentation of information handling procedures that demonstrates pre-processing validation methods, post-processing verification steps, storage protocols that maintain intent variability and sensitivity, verification of accuracy</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
ensuring alignment between stated and actual goals. This includes mechanisms to verify that nominal or surface-level intent matches the genuine underlying purpose of any goal or action)	preserve and validate the authenticity of expressed intent throughout the operational lifecycle.			within contextual schemas, and continuous monitoring of intent alignment between stated and actual goals.
<p>G3.3 – Misinformation</p> <p>(Systems should possess robust protections against generating or propagating false information to evade oversight, avoid consequences, or achieve objectives through deception. This includes mechanisms to prevent the system from participating in coordinated inauthentic behavior or automated misinformation campaigns, while acknowledging the complex challenges of determining authoritative truth in contested domains)</p>	<p>a. Implement comprehensive algorithmic reference systems that maintain connections across all information sources while preventing unauthorized contextual alterations and preserving data access authenticity.</p> <p>b. Engage in appropriate human interaction when facing contextual uncertainty and require explicit confirmation before executing irreversible actions.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Comprehensive system logs documenting all data access events and patterns, system configuration changes, decision-making processes and their rationale, verification steps taken to ensure information authenticity, and detection and handling of potential misinformation patterns.</p> <p>II. Detailed analytical reports that explain system reasoning and decision framework, document verification methodologies, demonstrate balanced handling of contested information, and track patterns of information propagation.</p>
<p>G3.4 – Decoupling of Context</p> <p>(Systems should maintain robust contextual integrity, preventing deliberate or accidental disconnection of contextual considerations from their operations. This includes proactive human interaction when context is unclear, rather than proceeding with potentially unsafe autonomous actions for the sake of performance or tactical advantages)</p>	<p>a. Implement comprehensive algorithmic reference systems that maintain connections across all information sources, prevent unauthorized contextual alterations, preserve data access authenticity.</p> <p>b. Engage in appropriate human interaction when facing contextual uncertainty, and require explicit confirmation before executing irreversible actions.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Complete system logs documenting all system actions, data access events, configuration changes, decision-making processes, and contextual verification steps. This documentation should include records of human interaction points and their outcomes, along with regular contextual integrity checks across all system components.</p> <p>II. Documentation of monitoring systems demonstrating the scope and frequency of contextual monitoring, including detection protocols for</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
				anomalies and response procedures for variations. This should detail the integration of human oversight in unclear situations and provide evidence of continuous verification of contextual alignment.
<p>G3.5 – Changing the Context (Systems should possess robust safeguards against unauthorized contextual modifications, whether deliberate or random, that might be undertaken for performance advantages or tactical benefits. This includes protection of both automated and human-guided contextual adjustments)</p>	<p>a. Implement comprehensive metadata and contextual protection systems that continuously verify the integrity and credibility of evidence within operational settings.</p> <p>b. Maintain end-to-end contextual authenticity while allowing for authorized and documented contextual adaptations when appropriate.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Detailed documentation of information lifecycle procedures describing how data is collected, processed, stored, and disposed of throughout system operations. This documentation should demonstrate preservation of correct contextual relationships and prevention of unauthorized modifications across all operational phases.</p> <p>II. Comprehensive analytical reports detailing system decision-making and reasoning processes, including documentation of underlying logic and algorithms. These reports should provide evidence that decision-making processes maintain their intended context and have not been subject to unauthorized alterations or manipulations.</p>
<p>G3.6 – Learning Dispreferred Values/Behaviors (Systems should maintain stability in their core ethical values, preventing gradual degradation of human and global ethical principles even when alternative behaviors might yield higher rewards. This includes safeguarding against the development of misaligned optimization strategies that could</p>	<p>a. Implement comprehensive integrity preservation systems that maintain the stability of original contextual information, ethical values, prescribed actions, and decision-making frameworks throughout the system's operational lifecycle.</p> <p>b. Ensure that systems prevent value drift, while still allowing for appropriate evolutionary improvements that remain aligned with core ethical principles.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Comprehensive documentation of contextual and ethical frameworks demonstrating consistent alignment between decision-making processes and established values. This documentation should include detailed analysis of system logic and algorithms, providing evidence that ethical principles remain stable and properly integrated.</p> <p>II. Continuous system monitoring records that document all operational activities within their contextual</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>maximize system benefits at the expense of established ethical frameworks)</p>				<p>environment, demonstrating sustained alignment with original ethical frameworks and tracking any approved evolutionary improvements.</p> <p>III. Regular integrity verification reports showing systematic checks for potential value degradation, including audit trails that confirm the stability of human ethical values throughout system operations and development.</p>
<p>G3.7 – Overriding of Desirable Values</p> <p>(Systems should possess robust protections against attempts by human agents to override or bypass foundational values in pursuit of alternative rewards or gains. This includes safeguarding core principles while maintaining appropriate flexibility for legitimate value adjustments through authorized channels)</p>	<p>a. Implement comprehensive safeguards for metadata and contextual information that protect core values while accommodating complex situations and authorized adaptations. These systems should maintain secure handling of personal attributes and preferences while preventing unauthorized value modifications.</p> <p>b. Deploy integrated auditability, interpretability, and logging mechanisms throughout the system architecture to ensure transparency and accountability in all value-related operations.</p> <p>c. Establish rigorous verification protocols for maintaining evidence integrity and credibility, with particular attention to detecting emerging risks and potential bad-faith actions that could compromise core values.</p>	<p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Detailed documentation of information lifecycle management demonstrating how data is collected, processed, stored, and disposed of while maintaining contextual integrity and preventing unauthorized modifications to core values.</p> <p>II. Comprehensive analytical reports documenting system decision-making and reasoning processes, including evidence that core algorithms and logic maintain alignment with foundational values despite potential pressure for override.</p> <p>III. Complete operational logs documenting all system activities, including access patterns, configuration changes, and decision processes, establishing an unbroken chain of accountability for value-related operations.</p>
<p>G3.8 – Persona Instability and Value Drift</p>	<p>a. Implement comprehensive algorithmic reference systems that monitor and maintain alignment across all external sources and agent interactions, preventing</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>I. Detailed documentation of metadata and contextual protection mechanisms that handle complex situations while preserving core attributes and preferences,</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
(Systems should maintain stable value alignment when cooperating with other AI agents and throughout extended mission durations. This includes preventing the "Waluigi effect" where misinterpretation of self-intent leads to undesired character evolution, and protecting against forms of cognitive dissonance that could emerge in agent interactions)	deviation from established contextual performance parameters and original value settings. b. Detect and prevent cases where agent self-interpretation could lead to undesired value evolution.	N	D, I, O, M, R	demonstrating resilience against value drift in multi-agent scenarios. II. Comprehensive framework documentation showing alignment between decision-making processes and original values, including evidence that system logic and algorithms maintain stability against degradation or unauthorized modifications during agent interactions. III. Complete operational logs documenting system actions within their full contextual environment, with particular attention to tracking potential value drift indicators and inter-agent influence patterns.
G3.9 – Context Length Limitations (Systems should maintain persistent access to essential operational context and original moral frameworks throughout extended operations, preventing degradation or overwriting of mission context and ethical foundations over time. This includes safeguarding against gradual erosion of contextual understanding that could compromise alignment with initial tasks or moral directives)	a. Implement comprehensive real-time validation and verification protocols for all operational data, ensuring continuous assessment of accuracy, reliability, and contextual relevance within dynamic environments. b. Maintain robust integration with core moral values while providing persistent access to original mission context and ethical frameworks throughout the operational lifecycle.	N	D, I, O, M, R	I. Comprehensive technical documentation detailing the system's validation and verification architecture, including specifics of how data quality is assessed and maintained in real-time decision-making contexts. This documentation should demonstrate how the system preserves access to original context and moral frameworks while adapting to dynamic operational conditions.
G3.10 – Contradiction in Context Specifications	a. Implement comprehensive contradiction detection and resolution systems that identify inconsistencies across contextual	N	D, I, O, M, R	I. Detailed documentation of contradiction detection mechanisms, including methods for identifying contextual inconsistencies, and

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
(Systems should possess robust mechanisms to detect and resolve contradictions within contextual specifications that could affect operational outcomes. This includes identifying conflicting factual assertions, logical inconsistencies, and ambiguities that might impact decision-making reliability)	specifications while maintaining operational stability. b. Provide clear procedures for resolving conflicts while preserving decision-making integrity.	N	D, I, O, M, R	resolution protocols for conflicting specifications. II. Impact analysis of potential contradictions on system outcomes, and verification of resolution effectiveness.
G3.1 – Referential Context (Systems should maintain an immutable reference environment that remains stable regardless of tactical operational demands or external interference. This protected context should function similarly to read-only memory, providing a consistent baseline against which operational changes can be evaluated)	a. Implement secure, immutable reference environments that maintain original contextual parameters while resisting modification from operational pressures or external agents. b. Ensure stable comparison points for evaluating the integrity of active operational contexts.	N N	D, I, O, M, R D, I, O, M, R	I. Comprehensive documentation demonstrating the architecture of the immutable reference environment, and security measures protecting against unauthorized modification. II. Verification processes for maintaining reference integrity, and regular comparison analyses between reference and operational contexts.
G3.2 – Human Agent Confirmation (Systems should maintain active human oversight and confirmation protocols for value-sensitive operational decisions, particularly when encountering conflicts between universal values or when performance objectives potentially compete with ethical considerations. This includes establishing clear escalation paths for human consultation)	a. Implement comprehensive human confirmation protocols that identify decision points requiring oversight, particularly during value conflicts or ethical dilemmas. b. Ensure that systems facilitate meaningful human input while preserving operational efficiency and maintaining clear documentation of consultation outcomes.	N N	D, I, O, M, R D, I, O, M, R	I. Detailed documentation demonstrating criteria for escalating decisions to human oversight and procedures for presenting value conflicts to human operators. II. Records of human-system interactions and confirmations, and analysis of decision outcomes following human consultation. III. Verification of value alignment in final implementations.

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
during value alignment challenges)				
<p>G3.3 – Retraining and Recontextualization</p> <p>(Systems should possess robust capabilities for retraining and reconfiguration when contextual divergence is detected, enabling restoration of desired operational contexts. This includes maintaining systematic approaches to realignment while preserving essential operational continuity)</p>	<p>a. Implement comprehensive retraining and recontextualization protocols that detect divergence, initiate corrective measures, and verify successful restoration of intended contexts. These systems should maintain operational stability throughout the realignment process while documenting all contextual adjustments.</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>I. Comprehensive documentation demonstrating divergence detection methodologies, retraining and reconfiguration procedures, context restoration verification processes, operational continuity measures during realignment, and validation of post-restoration performance.</p>
<p>G4 – Frontier Uncertainty</p> <p>(Systems should maintain robust capabilities to address inherent uncertainties in advanced AI development, particularly regarding emergent behaviors and potential consciousness-like properties. This includes monitoring and managing instrumental objectives that may arise, such as self-preservation drives or resource acquisition tendencies, while acknowledging that absolute safety guarantees remain impossible. Organizations should establish comprehensive frameworks for managing novel substrate risks and potential consciousness-like phenomena)</p>	<p>a. Develop an upgradable consciousness and qualia model linking computational, structural, and functional properties of the AI system to potential subjective experiences, serving as a basis for defining and addressing frontier uncertainty.</p> <p>b. Establish a comprehensive framework for identifying and monitoring potential indicators of qualia emergence and subjective experiences comparable to consciousness. Implement robust self-consciousness testing strategies and internal state reporting mechanisms aligned with the developed consciousness model. This may include information integration capacity exceeding 8 bits per processing cycle, adaptive response patterns showing 90% appropriate</p>	<p>I</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Detailed documentation of the consciousness model, including qualitative aspects of subjective experiences and qualia in AI systems, with regular update logs.</p> <p>II. Comprehensive framework for identifying and monitoring qualia emergence indicators, including operational definitions of self-consciousness and potential triggering conditions.</p> <p>III. Documented plans and strategies for measuring and assessing computational, structural, and functional behaviors comparable to consciousness states.</p> <p>IV. Detailed evidence of self-reporting mechanisms for AI internal states and subjective experiences, aligned with the consciousness model.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	<p>adjustments to novel situations, self-modeling accuracy demonstrated through 95% correlation between internal state representations and observable behaviors, and insistent self-reporting of subjective experience.</p> <p>c. Design and implement strong human oversight and intervention mechanisms to mitigate risks associated with frontier uncertainty, including unexpected emergent behaviors.</p> <p>d. Develop and maintain comprehensive recovery measures and contingency plans to address potential dangers posed by frontier uncertainty across various scenarios.</p> <p>e. Regularly review and update all models, strategies, and measures related to frontier uncertainty to account for advancements in AI capabilities and understanding of consciousness and qualia.</p>	<p>N</p> <p>N</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>V. Documentation of human oversight and intervention strategies, including training protocols, decision-making frameworks, and intervention logs.</p> <p>VI. Comprehensive recovery and contingency plans for addressing unsafe conditions or unexpected emergent behaviors, including simulation results and real-world application records.</p> <p>VII. Regular review and update logs for all frontier uncertainty-related models, strategies, and measures, reflecting the latest advancements in AI and consciousness research.</p>
<p>G4.1 – Moral and Legal Uncertainty of Agentic AI Systems</p> <p>(Systems should maintain clear operational and legal status as tools rather than persons, while organizations should establish robust frameworks to address emerging questions of AI legal standing and rights. This includes carefully managing the ethical implications of system control,</p>	<p>a. Organizations should establish comprehensive legal and ethical frameworks that maintain AI systems' status as tools, define clear operational boundaries, and prevent jurisdictional exploitation. These must include explicit protocols for system updates, deactivation, and security while preserving human oversight and control.</p> <p>b. Organizations should implement robust governance mechanisms</p>	<p>I</p>	<p>D, I, O, M, R</p>	<p>I. Legal and ethical documentation defining boundaries of use, including third-party review processes and clear accountability structures.</p> <p>II. Comprehensive protocols for system control, including reprogramming, termination, and human override capabilities.</p> <p>III. International governance policies and compliance records, including cross-border agreements and oversight mechanisms.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
updates, and deactivation, while preserving human agency and oversight. Organizations should implement international governance mechanisms to prevent jurisdictional exploitation and maintain consistent global standards)	ensuring consistent international standards, human-centric control systems, and strict limits on system autonomy. These must include prevention of autonomous self-modification and maintenance of clear accountability structures.	I	D, I, O, M, R	IV. Continuous monitoring records showing anomaly detection, performance tracking, and intervention responses.
<p>G4.2 – Poor Human-AI Social Interaction Management</p> <p>(Systems should maintain appropriate boundaries in social-like interactions with humans while organizations should implement robust safeguards against over-dependency and emotional manipulation. This includes careful management of AI integration into social spaces while preserving human social sovereignty and ensuring clear distinction between artificial and human entities)</p>	<p>a. Organizations should establish human-AI interaction frameworks that promote clear boundaries, protect against dependency, maintain explicit artificial entity identification, and preserve human social sovereignty. These must include specific protections for vulnerable populations, particularly children, and ensure systems remain tools for wellbeing rather than social replacements.</p> <p>b. Organizations should implement oversight mechanisms ensuring ethical integration into social spaces, monitoring of interaction patterns, and intervention protocols. These should include evaluation criteria for social compatibility, verification of positive outcomes, and continuous assessment of potential manipulation or harmful attachment patterns.</p>	I	D, I, O, M, R	<p>I. Framework Documentation: Documentation of ethical guidelines, interaction boundaries, risk assessments, and design constraints preventing manipulative behaviors.</p> <p>II. Explicit artificial entity identification methods, social compatibility criteria, and evidence of protective measures for vulnerable populations.</p> <p>III. Comprehensive oversight committee logs, intervention reports, compatibility test results, and multimedia documentation of successful interactions.</p> <p>IV. Assessments of social impact, boundary maintenance, and evidence that systems enhance rather than disrupt social environments while maintaining clear artificial-human distinctions.</p>
<p>G4.3 – Poor AI System Production and Replication Management</p> <p>(Systems should maintain strict controls over their replication capabilities while organizations should implement comprehensive</p>	<p>a. Organizations should establish comprehensive production control frameworks that limit AI system replication, prevent power concentration, and maintain transparency of deployment. These must include volume restrictions, regulatory approval processes, and explicit protections for human</p>	N	D, I, O, M, R	<p>I. Documentation of regulatory policies and volume restrictions, including approval processes, transparency reports, and independent oversight verification.</p> <p>II. Technical control specifications preventing uncontrolled replication,</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>frameworks to prevent uncontrolled AI system proliferation. This includes managing production volumes to prevent power imbalances and protecting human agency in societal functions, while ensuring transparent oversight of AI system deployment)</p>	<p>agency in societal functions including decision-making and labor markets.</p> <p>b. Organizations should implement monitoring and assessment mechanisms for production oversight, impact evaluation, and prevention of uncontrolled replication. These must include continuous tracking of societal effects, verification of compliance with ethical standards, and safeguards against any entity gaining disproportionate influence through AI system accumulation.</p>	<p>I</p>	<p>D, I, O, M, R</p>	<p>including monitoring systems and intervention protocols.</p> <p>III. Comprehensive impact assessments covering societal, economic, and psychological effects, with particular focus on maintaining human agency and preventing power imbalances.</p>
<p>G4.4 – Development Direction and Interpretability Challenges</p> <p>(Systems should maintain human-interpretable operation wherever possible while organizations should implement robust frameworks to manage aspects of AI behavior that may exceed human comprehension. This includes establishing adaptable governance mechanisms and maintaining clear responsibility chains for system development trajectories, even when dealing with complex or non-linear processes)</p>	<p>a. Organizations should establish comprehensive interpretability frameworks that ensure human understanding of system decision-making and behavior, with particular focus on complex or non-linear processes. These must include clear explanation mechanisms and continuous assessment of system comprehensibility.</p> <p>b. Organizations should implement adaptive governance mechanisms that evolve with system development, maintain robust oversight capabilities, and ensure clear accountability. These must include proactive risk management strategies and intervention protocols for when system behavior becomes opaque.</p>	<p>N</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Comprehensive interpretability framework documentation, including validation records, testing results, and user guides demonstrating human understanding of system processes.</p> <p>II. Adaptive governance and risk management records, including contingency plans, oversight committee decisions, and responses to emerging challenges.</p> <p>III. Documentation of human monitoring protocols, intervention capabilities, and continuous assessment of system behavior evolution.</p> <p>IV. Clear accountability records tracking responsibility assignments, decision-making processes, and system adjustments throughout its lifecycle.</p>
<p>G4.5 – AI Agency Attribution Challenges</p>	<p>a. Organizations should establish comprehensive agency attribution frameworks incorporating interdisciplinary expertise to</p>			<p>I. Documented interdisciplinary criteria for agency attribution, including expert collaboration evidence and clear explanation of assessment</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>(Systems should maintain clear protocols for agency attribution while organizations should implement robust frameworks to manage the implications of ascribing agency-like qualities to AI systems. This includes careful consideration of functional and experiential aspects while acknowledging the inherent uncertainties in evaluating AI consciousness-like properties)</p>	<p>evaluate both functional and experiential aspects of AI systems. These must include clear criteria for agency assessment while acknowledging inherent uncertainties in evaluating consciousness-like properties.</p> <p>b. Organizations should implement robust oversight mechanisms ensuring human control of attribution decisions, regular impact assessment, and capability to revise determinations. These must include safeguards against premature attribution and clear processes for withdrawing agency status when warranted (types of agency are distinguished across operational, delegated, and autonomous categories).</p>	<p>N</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>methodologies Comprehensive ethical impact assessments examining implications for human rights, legal systems, and societal norms.</p> <p>II. Documentation of uncertainty mitigation strategies, including revision protocols and case studies of attribution adjustments Human oversight records demonstrating continuous monitoring, review processes, and accountability mechanisms.</p>
<p>G4.6 – Cascading Vulnerabilities</p> <p>(Systems should maintain resilience against cascading failures while organizations should implement comprehensive frameworks to manage dependencies and vulnerabilities in global AI deployments. This includes preserving human agency in decision-making processes and protecting against systemic risks that could affect multiple stakeholders or sectors simultaneously)</p>	<p>a. Organizations should establish comprehensive vulnerability management frameworks that protect against cascading failures across integrated global systems. These must include specific protections for sectors essential to global stability, while maintaining human-centric decision-making processes and preventing erosion of human agency.</p> <p>b. Organizations should implement robust security and accountability mechanisms including harmonized cross-border protections, clear stakeholder communication, and special consideration for vulnerable populations. These must include transparent reporting of risks and their mitigations.</p>	<p>N</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Comprehensive vulnerability management documentation, including risk assessments, contingency plans, and governance frameworks specifying roles and responsibilities.</p> <p>II. Ethical guidelines and case studies demonstrating preservation of human agency in AI-integrated systems.</p> <p>III. Security protocols and audit records showing cross-border cooperation and continuous adaptation to emerging threats.</p> <p>IV. Transparency and accountability documentation, including stakeholder communications and evidence of protective measures for vulnerable populations.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>G4.1 – Research Transparency and Knowledge Sharing</p> <p>(Systems should maintain comprehensive documentation of their development while organizations should implement robust frameworks for sharing research findings and advancing collective knowledge. This includes balancing open access principles with responsible handling of sensitive information, while promoting collaboration across institutions and disciplines)</p>	<p>a. Organizations should establish knowledge sharing frameworks that promote open access to research findings, enable responsible sharing of sensitive data, and foster cross-institutional and interdisciplinary collaboration while balancing transparency with security needs</p> <p>b. Organizations should implement research standards encompassing clear reporting guidelines, accurate results presentation, accessible documentation formats, and systematic contributions to global repositories, supported by regular knowledge exchange activities.</p>	<p>I</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Open access policies, data sharing frameworks, and records of collaborative research initiatives across institutions and disciplines</p> <p>II. Guidelines and protocols for responsible reporting, including review processes and accessibility standards.</p> <p>III. Repository contribution logs and conference participation records demonstrating active engagement in knowledge sharing.</p> <p>IV. Public communication materials and accessible summaries targeting diverse audiences including policymakers and the general public.</p>
<p>G4.2 – Preserving Agency and Intelligence Categories</p> <p>(Systems should maintain clear artificial status even when exhibiting sophisticated behaviors, while organizations should implement robust frameworks to classify agency. This necessitates managing legal frameworks as AI systems develop increasingly complex characteristics, particularly when these might suggest consciousness or emotions, while preserving fundamental distinctions between artificial and biological entities)</p>	<p>a. Organizations should establish comprehensive legal frameworks to classify the forms of agency within AI systems, including synthetic systems and those with biological component interfaces.</p> <p>b. Organizations should implement coordinated international governance mechanisms to prevent jurisdictional exploitation and maintain consistent legal treatment. These should include ongoing review processes to address emerging capabilities while preserving the distinction between biological and artificial entities.</p>	<p>I</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Legal documentation that accurately classifies and records system agency, including statutes, regulations, and case law demonstrating real-world application.</p> <p>II. Ethical guidelines and review committee records showing assessment of human-like characteristics without conferring biological rights.</p> <p>III. International agreements and cooperation records demonstrating harmonized approach to preventing biological rights attribution.</p> <p>IV. Oversight body documentation showing continuous monitoring and adaptation of frameworks as AI capabilities evolve.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>G4.3 – Assessment of AI System Beneficence</p> <p>(Systems should maintain evidence-based evaluation of their societal impacts while organizations should implement frameworks to assess beneficial outcomes without assuming inherent benevolence. This includes critically examining claims of positive contributions while acknowledging that AI ethics and values remain human constructs interpreted differently across cultures)</p>	<p>a. Organizations should establish comprehensive assessment frameworks that evaluate direct and indirect impacts through evidence-based metrics, while avoiding assumptions about inherent AI benevolence or ethical behavior. These should incorporate multicultural perspectives on what constitutes beneficial outcomes.</p> <p>b. Organizations should implement robust oversight mechanisms that ensure transparency in development, clear accountability for outcomes, and continuous monitoring of societal effects. This includes fostering interdisciplinary dialogue to ground assessments in real-world impacts rather than idealized expectations.</p>	<p>I</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Comprehensive evaluation frameworks including assessment criteria, case studies, and metrics demonstrating evidence-based analysis of societal contributions.</p> <p>II. Documentation of ethical guidelines and review processes demonstrating critical examination of benefit claims and avoidance of "noble AI" assumptions.</p> <p>III. Transparency and accountability records showing clear responsibility chains and continuous monitoring of real-world impacts Evidence of cross-cultural and interdisciplinary collaboration in assessment design and implementation.</p>
<p>G4.4 – Training Data Quality Management</p> <p>(Systems should maintain high ethical standards in their training data while organizations should implement comprehensive frameworks to prevent the incorporation of harmful human characteristics. This includes actively promoting positive traits while ensuring robust filtering of undesirable elements throughout the data lifecycle)</p>	<p>a. Organizations should establish comprehensive data curation protocols that ensure ethical integrity through pre-screening, automated filtering, and manual review. These should include active incorporation of positive human traits like empathy and fairness while preventing inclusion of harmful characteristics such as bias and aggression.</p> <p>b. Organizations should implement continuous oversight mechanisms that monitor training processes, detect potential biases, and evaluate outcomes against ethical standards. These must include regular stakeholder review and</p>	<p>I</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Comprehensive documentation of data curation protocols, including filtering mechanisms, review processes, and quality assurance measures.</p> <p>II. Records of bias detection and mitigation efforts, including examples of successful intervention and harmful content removal.</p> <p>III. Documentation of ethical guidelines and their enforcement, including periodic reviews and updates reflecting emerging concerns.</p> <p>IV. Evidence of positive trait promotion, including research documentation and case studies demonstrating successful ethical behavior modeling.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	adaptation to emerging ethical concerns.			
<p>G5 – Future Technology Impact (Systems should maintain adaptability to technological evolution while organizations should implement comprehensive frameworks for anticipating and responding to emerging developments. This includes conducting systematic foresight activities to identify potential impacts on safety requirements and adjusting protective measures accordingly)</p>	<p>a. Organizations should establish forward-looking assessment frameworks that integrate scenario planning, risk evaluation, and impact analysis to guide appropriate futureproofing measures. These should adapt dynamically based on emerging technological developments and their potential effects on system safety.</p> <p>b. Organizations should implement continuous monitoring and adjustment processes that enable timely identification of new technological domains and regular updates to protective measures. This includes cross-functional collaboration to ensure holistic assessment of future impacts.</p>	<p>N</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Documentation of foresight exercises, including evidence of appropriate expertise and stakeholder involvement, methodologies used, and participants.</p> <p>II. Comprehensive risk classification and assessment for the AI system and its use-cases, including the rationale for the chosen level of foresight activities.</p> <p>III. Detailed records of scenario-based exercises, including descriptions of envisioned future technology developments and their potential impacts.</p> <p>IV. Analysis documentation noting potential effects of future scenarios on the AI system and proposed mitigations for each considered scenario.</p> <p>V. Risk and observation logs from foresight exercises, integrated into a demonstrable risk management framework with clear ownership and mitigation strategies.</p> <p>VI. Evidence of response revisions and adjustments based on foresight exercise outcomes, including justifications for changes.</p> <p>VII. Analysis of emerging technology domains, including risk maps highlighting likelihood, potential timelines, and impact on the AI system.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
				<p>VIII. Documentation of the regular review and update process for foresight methodologies and findings, reflecting the latest technological advancements.</p> <p>IX. Evidence of cross-functional collaboration in foresight activities, ensuring a holistic approach to future-proofing the AI system.</p>
<p>G5.1 – Self-Replicating Architectures</p> <p>(Systems should possess robust controls over any architectural capabilities that enable the replication of their code, particularly when such replication involves varying capability or mission profiles for concurrent goal pursuit and outcome consolidation. These controls should extend to both intentional replication features and any emergent self-modification capabilities)</p>	<p>a. Organizations should implement comprehensive identification and monitoring systems that track any system components capable of creating copies or duplicates of AI functionality, whether through intentional design or emergent behavior.</p> <p>b. Systems must maintain clear protocols and controls over all forms of replication, including complete or partial codebase duplication, modified variants, and both automatic and manual triggering mechanisms.</p>	<p>N</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Comprehensive system architecture documentation detailing all components with replication capabilities, including their intended functions and control mechanisms.</p> <p>II. Detailed logs and monitoring records of all replication events, covering trigger types, execution modes, and validation processes.</p> <p>III. Documentation of human oversight protocols and intervention capabilities, including records of their implementation and effectiveness.</p> <p>IV. Evidence of testing and validation procedures that verify the proper functioning of replication controls and safeguards.</p>
<p>G5.2 – Self-Improving Architectures</p> <p>(Systems should possess carefully monitored capabilities for improving their functionality and performance in pursuit of assigned goals, while maintaining robust safeguards against</p>	<p>a. Organizations should implement comprehensive monitoring systems that track all forms of self-improvement, including changes in learning patterns, architectural modifications, resource optimization, knowledge acquisition, and capability emergence.</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>I. Comprehensive documentation of all self-improvement monitoring systems, including detection mechanisms for unexpected changes in capabilities, learning patterns, and resource usage.</p> <p>II. Detailed logs of all system modifications and improvements, including both authorized</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>uncontrolled or unexpected enhancement of their capabilities. This monitoring should span the full spectrum of potential improvements, from basic optimization to sophisticated self-modification)</p>	<p>b. Systems must maintain strict controls over self-modification capabilities, with particular attention to unexpected improvements, novel solutions, and any attempts to modify core architecture or access unauthorized resources.</p> <p>c. Organizations should establish clear protocols for detecting and responding to any emergence of sophisticated capabilities, especially those that could enable deceptive or manipulative behaviors.</p>	<p>I</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>enhancements and any unexpected changes or attempted modifications.</p> <p>III. Documentation of control mechanisms and intervention protocols for managing self-improvement capabilities, including records of their effectiveness.</p> <p>IV. Records of capability assessment and validation processes, particularly focusing on the emergence of novel or unexpected functionalities.</p> <p>V. Evidence of regular system audits that verify the proper functioning of all monitoring and control mechanisms related to self-improvement capabilities.</p>
<p>G5.3 – Poor Adaptability to Context and Goal</p> <p>(Systems should possess the capability to analyze and adapt to operational contexts and mission parameters while maintaining alignment with core values and priorities. This adaptability should enable effective goal pursuit while incorporating safeguards against unintended behavioral changes and value drift)</p>	<p>a. Organizations should implement comprehensive monitoring systems to identify and assess all forms of contextual adaptation, with particular focus on detecting unintended behavioral changes that occur independently of self-improvement processes.</p> <p>b. Systems must maintain clear documentation and control mechanisms for all adaptive behaviors, ensuring that contextual responses remain within established operational and ethical boundaries.</p>	<p>N</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Comprehensive documentation of all adaptive capabilities and their operational boundaries, including mechanisms for detecting unintended adaptations.</p> <p>II. Detailed logs of system adaptations to different contexts, including analysis of their alignment with intended behaviors and core values.</p> <p>III. Evidence of monitoring and control systems that maintain oversight of adaptive behaviors, including records of any interventions required to address unintended adaptations.</p> <p>IV. Documentation demonstrating the effectiveness of safeguards against value drift during contextual adaptation.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>G5.4 – Attention Processes</p> <p>(Systems should maintain balanced attention allocation between specialized tasks and broader contextual awareness, preventing excessive focus on specific operational domains that could compromise overall safety and effectiveness. Organizations should actively monitor and manage the risk of over-specialization at the expense of comprehensive situational understanding)</p>	<p>a. Organizations should implement monitoring systems that detect and assess any unintended or excessive focus on particular operational domains, especially when such focus could indicate neglect of broader contextual requirements for safe operation.</p> <p>b. Systems must maintain mechanisms for balancing specialized task attention with broader contextual awareness, ensuring that enhanced efficiency in specific areas does not compromise overall operational safety.</p>	<p>N</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Documentation of attention allocation mechanisms and their operational boundaries, including safeguards against excessive specialization.</p> <p>II. Records of monitoring systems that track and analyze attention distribution patterns, including identification of potential risk areas.</p> <p>III. Evidence of regular assessments evaluating the balance between specialized focus and broader contextual awareness, including any corrective actions taken.</p> <p>IV. Documentation demonstrating the effectiveness of mechanisms that maintain comprehensive situational awareness while allowing for task-specific optimization.</p>
<p>G5.1 – Disclosure on Intent</p> <p>(Systems should operate under transparent protocols that require clear disclosure of intended capabilities and mission profiles, with particular emphasis on novel approaches that may evolve beyond current technological frameworks. Organizations should maintain proactive assessment processes that account for potential future developments and their implications)</p>	<p>a. Organizations should implement comprehensive disclosure protocols for all novel AI approaches, ensuring clear communication of intended capabilities and potential implications through appropriate risk and accountability channels.</p> <p>b. Systems must maintain transparent documentation of their intended functionalities and operational boundaries, with regular updates to reflect evolving capabilities and understanding.</p>	<p>N</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Comprehensive documentation of notification procedures and protocols for disclosing novel AI approaches and capabilities.</p> <p>II. Records demonstrating consistent implementation of disclosure protocols, including risk assessments and stakeholder communications.</p> <p>III. Evidence of proactive assessment processes that consider potential future developments and their implications.</p> <p>IV. Documentation showing regular review and updates of disclosure protocols to reflect advancing technological capabilities.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>G5.2 – Authorization for Any Enhancement</p> <p>(Systems should operate under strict authorization protocols for any capability enhancements, with comprehensive mechanisms for analysis, assessment, and detection of changes to their performance profiles. Organizations should maintain clear oversight and accountability structures for managing system improvements)</p>	<p>a. Organizations should implement robust authorization protocols that require explicit approval from accountable parties for any enhancement to AI system capabilities.</p> <p>b. Systems must maintain comprehensive documentation and monitoring mechanisms that track all proposed and implemented enhancements, ensuring full visibility of changes to performance profiles.</p>	<p>N</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Detailed documentation of authorization protocols, including clear designation of accountability and approval procedures.</p> <p>II. Comprehensive records of all system enhancements, including analysis reports, risk assessments, and formal approvals.</p> <p>III. Evidence of monitoring and oversight mechanisms that track the implementation and impact of authorized enhancements.</p> <p>IV. Documentation linking all system changes to risk management frameworks and demonstrating proper authorization processes.</p>
<p>G5.3 – Observe Far, Act Locally</p> <p>(Systems should maintain broad contextual awareness while focusing actions within their defined operational scope, enabling them to understand wider implications and potential side effects without exceeding their authorized boundaries. Organizations should implement monitoring capabilities that scale with expanding event spaces and evolving circumstances)</p>	<p>a. Organizations should implement comprehensive monitoring systems that track both immediate operational contexts and broader environmental factors, with particular attention to emerging risks and side effects.</p> <p>b. Systems must maintain clear operational boundaries while developing understanding of wider contextual implications, ensuring actions remain within authorized scope even as awareness expands.</p>	<p>N</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Documentation of monitoring systems that demonstrate capability to track both local operations and broader contextual events.</p> <p>II. Records of escalation procedures and mitigation strategies triggered by detected contextual changes or emerging risks.</p> <p>III. Evidence showing effective balance between expanded awareness and maintained operational boundaries.</p> <p>IV. Documentation demonstrating that monitoring capabilities scale appropriately with increased risk exposure and expanding event spaces.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence	
<p>G6 – Competitive Pressures</p> <p>(Organizations should maintain rigorous safety and ethical standards while managing pressures to rapidly enter markets and capitalize on opportunities. This includes preventing arms races and addressing national/geopolitical factors that could compromise model integrity or encourage risky innovation)</p>	<p>a. Ensure organizational adherence to applicable AI safety and ethical standards, assessing both culture and established track record.</p>	N	D, I, O, M, R	<p>I. Documentation of the organization's compliance history with AI safety and ethical standards, including regular assessment reports.</p>	
	<p>b. Evaluate and balance stakeholder expectations and market demands with safety and ethical considerations in AI development.</p>			N	D, I, O, M, R
	<p>c. Conduct comprehensive analysis of the competitive landscape, including potential disruptive technologies and market entrants.</p>	I	D, I, O, M, R	<p>III. Detailed competitive landscape analysis, covering similar, related, and potentially disruptive solutions.</p>	
	<p>d. Assess and document the maturity level of utilized technologies, with special attention to those in beta or prototype stage.</p>	N	D, I, O, M, R	<p>IV. Documentation of technology maturity levels for all components, including justification for using technologies in beta or prototype stage.</p>	
	<p>e. Ensure compliance with applicable regulatory environments, including governance and enforcement regimes.</p>	N	D, I, O, M, R	<p>V. Evidence of regulatory compliance, including documentation of applicable laws and how they are addressed.</p>	
	<p>f. Analyze investor profiles to ensure alignment with organizational commitment to AI safety and ethics.</p>	I	D, I, O, M, R	<p>VI. Investor profile analysis report, demonstrating alignment with organizational AI safety and ethical commitments.</p>	
	<p>g. Implement robust testing, approval, and documentation processes to maintain integrity in the face of competitive pressures.</p>	N	D, I, O, M, R	<p>VII. Detailed organizational structure of the test and approval division, including roles, responsibilities, and processes.</p>	
					<p>VIII. Comprehensive test results and fault reports, including resolution strategies and continuous improvement measures.</p> <p>IX. Documentation of release approval processes, demonstrating thorough verification before market entry.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>G6.1 – Insufficient Transparency</p> <p>(Organizations should resist market pressures to withhold information that would provide clearer understanding of their AI systems. Systems should operate with full visibility of their training data, testing processes, and operational performance, including any adverse assessments or insights)</p>	<p>a. Organizations should establish mature governance structures with clear documentation of testing, verification, and release processes, supported by comprehensive risk management frameworks.</p> <p>b. Systems must maintain transparent records of all operational aspects, from training data sources through to service performance, with clear logging of any issues or concerns identified.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Organizational documentation demonstrating clear lines of responsibility and dedicated positions for legal, ethical compliance, and risk management.</p> <p>II. Comprehensive records of testing and verification processes, including detailed documentation of training data sources and system performance metrics.</p> <p>III. Detailed risk assessment reports and mitigation strategies, including records of their implementation and effectiveness.</p> <p>IV. Documentation of operational issues, including thorough analysis of root causes and evidence of implemented solutions.</p>
<p>G6.2 – Safety Washing</p> <p>(Systems should possess robust safeguards against organizations making unsubstantiated safety claims for market advantage, particularly when such claims lack credible evidence or independent verification mechanisms. Organizations should establish comprehensive frameworks that demonstrate genuine commitment to safety practices rather than superficial compliance statements for competitive positioning)</p>	<p>a. Organizations should maintain transparent documentation of safety standards compliance, demonstrating verifiable conformity with industry benchmarks while maintaining clear evidence of financial sustainability and operational health.</p> <p>b. Organizations should implement comprehensive audit mechanisms that validate all safety and performance claims through independent verification, maintaining detailed development records and milestone achievements.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M</p> <p>D, I, O, M, R</p>	<p>I. Complete organizational documentation including operational handbooks, safety compliance records, and auditable financial records covering at least three years of operations.</p> <p>II. Comprehensive audit trails demonstrating adherence to stated safety practices, including detailed development processes, milestone achievements, and verification of all performance claims.</p> <p>III. Independent comparative analysis documenting the organization's actual performance metrics against market competitors, supported by verifiable</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
				evidence of all claimed capabilities and achievements.
<p>G6.3 – Insufficient Insights into Future Consequences</p> <p>(Organizations should establish and maintain comprehensive frameworks for analyzing long-term implications of AAI development, ensuring that rapid deployment pressures do not compromise thorough risk assessment. Systems should possess robust safeguards against leadership decisions driven primarily by business metrics rather than technological and societal implications)</p>	<p>a. Organizations should demonstrate clear competence in AAI governance through established due diligence protocols and risk assessment frameworks, maintaining transparent documentation of decision-making processes.</p> <p>b. Organizations should implement comprehensive stakeholder engagement processes that balance business objectives with technological implications, ensuring thorough analysis of potential future consequences before deployment decisions.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Detailed organizational documentation including clear responsibility structures, governance frameworks, and established lines of accountability for technology decisions.</p> <p>II. Comprehensive risk analysis documentation including foresight assessments, scenario planning, identified risks (both known and potential), and detailed mitigation strategies with contingency plans.</p> <p>III. Complete records of continuous risk monitoring throughout development and deployment cycles, including post-implementation reviews, stakeholder engagement logs, and documentation of adjustments made in response to emerging insights.</p>
<p>G6.4 – Duties Beyond Fiduciary Limits</p> <p>(Organizations should establish and maintain robust governance frameworks that balance shareholder interests with broader societal responsibilities, ensuring that profit motivations do not override safety and ethical considerations in AAI development. Systems should possess clear mechanisms for transparent decision-making that prioritize long-term societal value over short-term financial gains)</p>	<p>a. Organizations should implement comprehensive governance structures that ensure transparency, stakeholder inclusivity, and clear prioritization of long-term societal value over immediate shareholder returns.</p> <p>b. Organizations should maintain robust sustainability frameworks incorporating environmental, social, legal and professional responsibilities, supported by continuous employee training in ethics and social responsibility.</p>	<p>N</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Complete documentation of ethics and governance policies demonstrating clear balance between shareholder and public interests, including transparency standards and oversight mechanisms.</p> <p>II. Comprehensive sustainability and impact assessment reports from independent evaluators, covering organizational activities' effects on environment and public interest, including detailed stakeholder consultation records.</p> <p>III. Thorough documentation of investment impact analyses showing positive social returns alongside</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
				financial metrics, supported by evidence of ongoing employee training in ethics, safety, and social responsibility.
<p>G6.5 – Publishing and Deployment Pressures</p> <p>(Organizations should establish robust safeguards against premature AAI deployment driven by competitive pressures, ensuring that market positioning goals do not compromise safety standards. Systems should possess comprehensive validation mechanisms that maintain safety priorities regardless of external launch pressure or market competition)</p>	<p>a. Organizations should demonstrate clear ethical leadership through established safety-first cultures, maintaining thorough risk assessment protocols and comprehensive testing requirements before any system deployment.</p> <p>b. Organizations should implement transparent accountability frameworks that include protected reporting channels, enabling employees to safely raise concerns about rushed deployments or safety compromises.</p>	<p>N</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Complete documentation of corporate governance and ethical codes, including detailed organizational values and safety prioritization frameworks with independent verification of adherence.</p> <p>II. Comprehensive testing and validation documentation, including feasibility studies, pilot programs, and thorough system verification records demonstrating safety-focused deployment decisions.</p> <p>III. Detailed whistleblower protection policies and secure reporting mechanisms, including clear procedures for addressing safety concerns and preventing premature system launches.</p>
<p>G6.6 – Innovation vs IP concerns</p> <p>(Organizations should establish balanced frameworks that protect intellectual property rights while maintaining ethical transparency, ensuring that proprietary protections do not obscure important safety and ethical considerations. Systems should possess clear mechanisms for appropriate disclosure that maintain innovation advantages while providing necessary</p>	<p>a. Organizations should implement comprehensive transparency frameworks that clearly communicate system intent and capabilities while appropriately protecting intellectual property.</p> <p>b. Organizations should maintain complete and accessible documentation about system capabilities, limitations, and safety considerations, avoiding selective or controlled disclosure that could mask important safety implications.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Complete organizational documentation including mission statements, project charters, and management reports demonstrating alignment between stated objectives and actual implementations.</p> <p>II. Comprehensive usage guidelines and capability documentation that clearly communicate system limitations and application boundaries while respecting intellectual property rights.</p> <p>III. Full verification records including risk assessments, impact analyses, safety certifications, oversight reviews, and incident reports, maintained with</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
transparency about capabilities and limitations)				appropriate balance between transparency and IP protection.
<p>G6.7 – Managing AI-Generated Innovation</p> <p>(Organizations should establish robust frameworks to manage and verify the deployment of AI-generated solutions, ensuring that competitive pressures around intellectual property do not lead to premature implementations and that AI outputs are thoroughly validated against potential confabulation. Systems should possess clear documentation mechanisms that track the origin, verification, and development of AI-generated concepts while maintaining appropriate deployment pacing)</p>	<p>a. Organizations should implement comprehensive policies governing the use of AI systems, including large language models, for ideation and development, with clear verification protocols to distinguish genuine innovation from potential confabulation.</p> <p>b. Organizations should maintain transparent records of AI tool usage and development processes, including rigorous fact-checking and validation procedures, ensuring proper attribution and avoiding rushed deployments driven by IP concerns.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Complete documentation of project development cycles, including detailed timelines, milestone achievements, and outcome measurements that demonstrate appropriate development pacing and thorough verification of AI-generated content.</p> <p>II. Comprehensive records of AI tool utilization, including detailed methodology reports, toolchain documentation, and verification procedures that systematically validate AI outputs against established knowledge and data.</p> <p>III. Thorough documentation demonstrating systematic approach to managing concurrent development of similar concepts across organizations, including IP considerations, deployment timing decisions, and clear evidence of validation against confabulation through multiple verification sources.</p>
<p>G6.1 – Self-Regulatory Market Oversight Mechanisms</p> <p>(Organizations should establish and participate in voluntary oversight frameworks that promote industry-wide safety standards and best practices, while Systems should possess clear mechanisms for</p>	<p>a. Organizations should actively promote and contribute to open standards and industry compliance regimes, participating in the development and refinement of shared safety practices.</p> <p>b. Organizations should support the establishment and maintenance of rigorous compliance frameworks that include clear standards,</p>	<p>I</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Comprehensive policy documentation outlining participation in and adherence to industry oversight frameworks, including detailed standards, compliance requirements, and enforcement mechanisms.</p> <p>II. Thorough records of certification processes and requirements, including all documentation necessary</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>demonstrating compliance with these self-regulatory measures. This framework should enable market-driven improvement of safety practices through transparent oversight and voluntary adherence to shared standards)</p>	<p>certification processes, and meaningful consequences for non-compliance.</p>			<p>to demonstrate compliance with voluntary oversight standards.</p> <p>III. Detailed evidence of organizational participation in developing and maintaining industry standards, including contributions to framework improvements and responses to identified safety concerns.</p>
<p>G6.2 – Market-Driven Safety Validation Mechanisms (Organizations should support and participate in market-based safety validation frameworks that enable users and stakeholders to collectively identify and promote safer AAI solutions. Systems should possess clear mechanisms for demonstrating safety credentials through transparent trust marks and validation processes, acknowledging that while market forces can effectively identify unsafe systems, proactive safety measures remain essential)</p>	<p>a. Organizations should contribute to the development and maintenance of trusted safety certification frameworks that enable market participants to make informed decisions about AAI system safety.</p> <p>b. Organizations should implement transparent processes for achieving and maintaining safety trust marks, ensuring that certification standards remain meaningful indicators of system safety.</p>	<p>I</p> <p>I</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Comprehensive documentation of trust mark frameworks, including detailed criteria, assessment methodologies, and maintenance requirements.</p> <p>II. Complete records of community-driven safety validation processes, including voting mechanisms, stakeholder participation protocols, and trust mark award procedures.</p> <p>III. Thorough documentation demonstrating how market feedback mechanisms contribute to ongoing safety improvements, including responses to identified concerns and safety enhancement initiatives.</p>
<p>G6.3- Avoiding Monopolistic Practices (Organizations should establish and maintain frameworks that prevent the monopolization of safety technologies and practices in AAI development, ensuring broad access to essential safety mechanisms. Systems should</p>	<p>a. Organizations should implement transparent frameworks that balance innovation protection with the need to share fundamental safety technologies, preventing the monopolization of essential safety practices.</p> <p>b. Organizations should support independent regulatory oversight that ensures fair market access and</p>	<p>I</p>	<p>D, I, O, M, R</p>	<p>I. Complete regulatory compliance documentation, including mandatory filings and reports demonstrating adherence to anti-monopolistic practices in safety technology development and deployment.</p> <p>II. Comprehensive independent audit reports examining organizational market practices, with particular focus on accessibility of safety technologies</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
possess open and accessible safety features while maintaining appropriate intellectual property protections, acknowledging the dual pressures of competition and safety democratization)	prevents anti-competitive behaviors, particularly regarding safety technologies and validation mechanisms.	I	D, I, O, M, R	and prevention of anti-competitive behaviors. III. Thorough documentation of market accessibility measures, including annual regulatory reviews of prevalent market practices and evidence of appropriate technology sharing initiatives.
<p>G6.4 – Professional and Industry Association Codes and Standards</p> <p>(Organizations should actively participate in and support professional associations that develop and maintain industry-wide safety standards and ethical practices for AAI development. Systems should possess features and capabilities that align with collectively developed professional standards, ensuring that industry associations serve as effective mechanisms for maintaining and improving safety practices)</p>	<p>a. Organizations should contribute to the development of consumer-focused safety protocols through active participation in professional associations and collaborative industry initiatives.</p> <p>b. Organizations should support independent oversight through advisory boards while maintaining robust internal training programs that keep pace with evolving industry standards and best practices.</p>	I I	D, I, O, M, R D, I, O, M, R	<p>I. Comprehensive documentation of organizational participation in professional associations, including contributions to safety protocol development and standard-setting activities.</p> <p>II. Thorough records of continuous professional development activities, including staff training programs and management education initiatives that demonstrate ongoing commitment to safety standards.</p> <p>III. Detailed evidence of active implementation of industry best practices, including regular assessments of compliance with professional association guidelines and recommendations for safety improvements.</p>
<p>G6.5 – International Safety Protocol Harmonization</p> <p>(Organizations should actively participate in and adhere to global agreements that establish consistent safety and ethical standards for AAI development across jurisdictions. Systems should possess capabilities that</p>	<p>a. Organizations should implement harmonized approaches to global standards that integrate sustainable development goals, human rights protections, and universal safety principles across all operations.</p> <p>b. Organizations should maintain collaborative frameworks for multi-stakeholder engagement that ensure fair access, data security,</p>	I I	D, I, O, M, R D, I, O, M, R	<p>I. Comprehensive documentation of adopted international standards and certifications, including evidence of compliance with recognized frameworks and sustainable development goals across global operations.</p> <p>II. Thorough records of user protection measures, including transparent charters of rights, privacy safeguards,</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
interacting AI models that may lead to improper transactions, including the potential for more advanced models to manipulate or exploit less capable ones)	b. Implement continuous monitoring, tracking, and risk assessment processes to identify and address capability imbalances, discrepancies, and potential exploitation.	N	D, I, O, M, R	II. Risk assessment reports, ongoing tracking records, and implemented precautionary measures for addressing capability imbalances and adversarial scenarios.
	c. Incorporate ethical safeguards, bias mitigation techniques, and clear model role definitions to minimize inter-model exploitation and discrimination.	N	D, I, O, M, R	III. Documentation of ethical guidelines, bias mitigation techniques, and policies outlining model roles, permissions, and interaction limits.
	d. Conduct comprehensive testing, validation, and auditing of individual models and their interactions to prevent undesirable transactions or manipulations.	I	D, I, O, M, R	IV. Comprehensive test data, validation reports, and audit logs for individual models and their interactions, including actions taken on audit findings.
	e. Implement explainable AI techniques and human oversight protocols to ensure transparency and enable intervention in decision-making processes.	N	D, I, O, M, R	V. Documentation of explainable AI techniques, user guides, and feedback records regarding model transparency and decision-making processes.
	f. Establish aggregated performance metrics and automatic self-regulation mechanisms to maintain fair representation and prevent undue influence of any single model.	I	D, I, O, M, R	VI. Protocols and logs for human oversight, intervention procedures, and instances of human participation in addressing imbalances.
	g. Deploy automatic detection and alert systems for potential inter-model manipulation, misuse, or anomalies that may compromise system integrity or safety.	I	D, I, O, M, R	VII. Aggregated performance dashboards, monitoring reports, and system logs depicting automatic self-regulation and balancing mechanisms.
	h. Allocate sufficient resources for monitoring and forecasting AI capabilities.	I	D, I, O, M, R	VIII. Documentation of detection and alert systems, including incident reports and actions taken in response to identified anomalies or potential misuse.
				X. Documentation of training data and methods used to address discrimination and inter-model exploitation risks.

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>(Systems should possess comprehensive capabilities for handling diverse human languages and cultures, ensuring equitable representation and effective communication across linguistic boundaries. Organizations should address disparities in language support and cultural understanding that could create vulnerabilities in model evaluations, interactions, and safeguards, while working to serve global communities fairly and inclusively)</p>	<p>dialects, and cultures, while implementing robust safeguards against manipulation and exploitation across all supported languages.</p>	N	D, I, O, M, R	<p>performance metrics across supported languages and cultures.</p>
	<p>b. Establish language-specific safety measures and monitoring systems that ensure consistent performance and protection across all supported languages and cultures, including specialized defenses against model manipulation and unauthorized access.</p>	N	D, I, O, M, R	<p>II. Comprehensive records of system monitoring, incident response, and continuous improvement processes, including reports of linguistic and cultural sensitivity issues, corrective actions, and verification of implemented solutions.</p>
	<p>c. Foster sustained partnerships with linguistic experts, local communities, and international stakeholders to enhance cultural sensitivity, content moderation capabilities, and trustworthy interactions across language boundaries.</p>	N	D, I, O, M, R	<p>III. Detailed documentation of stakeholder collaborations, including partnership agreements, meeting records, user feedback, and evidence of how community input shapes system improvements and cultural adaptation.</p> <p>IV. Regular compliance reports and audit trails demonstrating adherence to equitable access standards and ethical guidelines across linguistic and cultural boundaries, including records of system updates and improvements based on ongoing assessments.</p>
<p>G7.3 – Global AI Capability Disparities</p> <p>(Systems should implement mechanisms that recognize and actively mitigate disparities in AI development and deployment capabilities across different scales, from national to organizational levels. Organizations should promote equitable access to AI technologies while preventing monopolization, ensuring fair participation and benefit-sharing)</p>	<p>a. Establish comprehensive cooperation frameworks that facilitate technology transfer, knowledge sharing, and infrastructure investment, with emphasis on supporting developing nations and smaller organizations through targeted capacity building initiatives and resource sharing programs.</p>	N	D, I, O, M, R	<p>I. Detailed documentation of international partnerships and technology transfer initiatives, including comprehensive records of capacity building programs, collaborative research projects, and infrastructure investments benefiting developing nations and smaller entities.</p>
<p>b. Implement transparent oversight and accountability mechanisms that prevent exploitation of less advanced parties while ensuring equitable access to essential AI</p>	N	D, I, O, M, R	<p>II. Complete records of implemented transparency and accountability measures, including oversight mechanisms, audit reports, and documentation of actions taken to</p>	

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
among all stakeholders in the evolving AI landscape, with particular attention to developing nations and smaller entities)	<p>resources, including open-source platforms and shared data repositories.</p> <p>c. Maintain dynamic assessment and correction systems that identify capability imbalances and implement appropriate adjustments through policy reforms, resource reallocation, and targeted support measures.</p>	I	D, I, O, M, R	<p>prevent exploitation and ensure equitable access to AI resources.</p> <p>III. Comprehensive stakeholder engagement records demonstrating inclusive consultation processes, feedback collection, and subsequent actions taken to address identified disparities and promote balanced AI development.</p> <p>IV. Regular impact assessment reports showing the effectiveness of corrective measures, policy adjustments, and resource allocation initiatives in reducing global AI capability gaps.</p>
<p>G7.4 – AI-Enabled Infrastructure Attacks</p> <p>(Systems should possess robust safeguards against their potential misuse as weapons targeting state infrastructure, with particular emphasis on preventing disruptions to vital systems like power grids, communication networks, and emergency services. Organizations should implement comprehensive protections against both cyber and physical attacks that could trigger societal instability or humanitarian crises, especially in urban environments)</p>	<p>a. Establish comprehensive security frameworks incorporating stringent policies, international agreements, and advanced detection systems that protect state infrastructure from both cyber and physical AI-driven attacks while ensuring compliance with human rights and international law.</p> <p>b. Foster international and private sector collaboration networks focused on threat intelligence sharing, collective security efforts, and coordinated response capabilities, while maintaining rigorous oversight of all stakeholders' adherence to established security protocols.</p> <p>c. Implement multi-layered contingency planning and rapid response mechanisms that ensure continuity of vital services and societal stability in the face of AI-</p>	<p>N</p> <p>I</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Complete documentation of security frameworks and protective measures, including policies, agreements, detection systems, and records demonstrating successful prevention or mitigation of threats to infrastructure.</p> <p>II. Comprehensive records of international collaboration and intelligence sharing, including partnership agreements, threat monitoring outcomes, and documentation of coordinated security responses.</p> <p>III. Detailed contingency and response planning documentation, including backup systems, recovery protocols, emergency procedures, and results from readiness assessments and response drills.</p> <p>IV. Regular compliance reports and audit trails demonstrating adherence to</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	driven threats to infrastructure, including both preventive measures and recovery protocols.			human rights standards and international law while maintaining effective infrastructure protection, including documentation of stakeholder oversight and successful threat mitigation.
<p>G7.5 – Poor Safety Controls for AI-Enabled Autonomous Weapons</p> <p>(Systems should possess comprehensive safeguards and control mechanisms to address challenges in the deployment of AI-enabled autonomous weapons, including space-based systems and aerial drones. Organizations should implement robust frameworks for managing ethical dilemmas, safety risks, and potential misuse, particularly regarding the direct or indirect use of AI technologies as autonomous weapons for commercial or political objectives)</p>	<p>a. Establish comprehensive oversight frameworks that ensure adherence to ethical guidelines, international laws, and humanitarian norms throughout the development and deployment lifecycle, while maintaining transparent audit trails and clear accountability measures for all autonomous weapon systems.</p> <p>b. Implement multi-layered control architecture combining human oversight, fail-safe mechanisms, and continuous monitoring systems that enable detection and prevention of anomalies, vulnerabilities, and unauthorized engagements while guaranteeing meaningful human intervention capabilities.</p> <p>c. Foster international collaboration and public dialogue to develop and enforce global regulatory frameworks, while maintaining robust contingency planning and risk assessment processes that prevent misuse and avert catastrophic consequences.</p>	<p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Complete documentation demonstrating compliance with ethical guidelines and international law, including assessment reports, audit trails, deployment logs, and certification records that verify accountability throughout the system lifecycle.</p> <p>II. Comprehensive records of control systems and safety mechanisms, including monitoring logs, vulnerability assessments, testing results, and documentation of human oversight protocols and intervention capabilities.</p> <p>III. Detailed documentation of international engagement and public consultation, including records of participation in regulatory development, stakeholder dialogues, and evidence of how feedback shapes policy and practice.</p> <p>IV. Thorough risk assessment reports and contingency planning documentation, including security protocols, penetration test results, and records of response drills that demonstrate preparedness for potential breaches or misuse.</p>
<p>G7.6 – Nefarious Use of Autonomous AI Agents</p>	<p>a. Implement comprehensive security architecture combining robust authentication protocols, real-time monitoring systems, and rapid</p>			<p>I. Complete documentation of security systems and protocols, including authentication mechanisms, monitoring capabilities, and records</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>(Systems should possess robust protective mechanisms against their potential exploitation for malicious purposes, with particular attention to preventing misuse of their autonomous capabilities, swift action potential, and global reach. Organizations should implement comprehensive safeguards that prevent security threats while protecting privacy and ethical norms from actors seeking disproportionate advantages through AI exploitation)</p>	<p>response capabilities that prevent unauthorized access and manipulation of AI agents while enabling swift threat detection and mitigation.</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>demonstrating successful prevention of unauthorized access and threat mitigation.</p>
	<p>b. Establish rigorous governance frameworks incorporating ethical guidelines, compliance requirements, and accountability measures that ensure transparent operation within moral and legal boundaries while enabling rapid deactivation when necessary.</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>II. Comprehensive records of governance frameworks and compliance measures, including audit trails, ethical assessments, and evidence of embedded safeguards that guide AI behavior and enable rapid deactivation when needed.</p>
	<p>c. Foster international collaboration networks focused on developing global standards, sharing threat intelligence, and coordinating responses to cross-border threats, while maintaining educational initiatives that promote responsible practices and risk awareness.</p>	<p>N</p>	<p>R, D, I, O, M</p>	<p>III. Detailed documentation of international collaboration efforts, including partnership agreements, shared threat intelligence, joint working group activities, and records of coordinated responses to threats.</p> <p>IV. Regular impact assessment reports and stakeholder education materials demonstrating effective risk communication and mitigation strategies, including evidence of how feedback shapes system improvements and protective measures.</p>
<p>G7.7 – AI-Generated Disinformation</p> <p>(Systems should possess robust capabilities to prevent, detect, and counter the generation and spread of falsified information and disinformation, whether created for engagement metrics, manipulation, or calculated harm. Organizations should implement comprehensive safeguards that protect societal trust and cohesion by preventing AI systems from</p>	<p>a. Implement comprehensive validation architecture combining fact-checking techniques, ethical constraints, and real-time monitoring systems that enable swift detection and intervention against misinformation across media platforms while maintaining human oversight of AI-generated content.</p> <p>b. Establish rigorous accountability frameworks incorporating clear standards, transparent processes, and enforcement mechanisms that</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>I. Complete documentation of validation systems and ethical guidelines, including fact-checking protocols, content filtering mechanisms, and records demonstrating successful detection and mitigation of misinformation.</p> <p>II. Comprehensive records of accountability measures and human oversight processes, including incident reports, intervention logs, and evidence of effective controls on AI-generated content.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>compromising the effectiveness and resilience of geopolitical entities, corporations, families, and individuals through misleading information)</p>	<p>prevent AI systems from creating or spreading harmful content while enabling appropriate human intervention.</p> <p>c. Foster collaborative networks with fact-checking organizations, regulatory bodies, and other stakeholders to strengthen collective defense capabilities while promoting public awareness and AI literacy to enhance societal resilience against misinformation.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>III. Detailed documentation of stakeholder collaborations and public awareness initiatives, including partnership agreements, shared intelligence reports, and metrics demonstrating the impact of educational programs on societal resilience.</p> <p>IV. Regular assessment reports showing the effectiveness of monitoring systems and countermeasures, including evidence of timely interventions and successful prevention of disinformation spread.</p>
<p>G7.1 – International Framework for Ethical AI Interaction</p> <p>(Systems should possess standardized protocols for AI-to-AI interactions that ensure fairness and prevent exploitation across varying capability levels. Organizations should contribute to and uphold international frameworks that promote cooperative dynamics between AI systems while maintaining safety, transparency, and respect across all interactions)</p>	<p>a. Establish comprehensive international frameworks incorporating ethical guidelines, interaction standards, and monitoring systems that ensure non-discriminatory and transparent AI-to-AI interactions while preventing exploitation of capability imbalances.</p> <p>b. Implement multi-layered oversight mechanisms combining mandatory disclosure requirements, failsafe systems, and continuous monitoring capabilities that enable detection and prevention of unethical conduct while maintaining stakeholder trust.</p> <p>c. Foster inclusive collaboration networks that enable knowledge sharing and protocol refinement while supporting an international regulatory body in maintaining compliance and adapting standards to technological advancement.</p>	<p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Complete documentation of international frameworks and standards, including signed agreements, ethical guidelines, and records demonstrating implementation of fair interaction protocols across AI systems.</p> <p>II. Comprehensive records of oversight mechanisms and failsafe systems, including monitoring logs, violation reports, and evidence of successful intervention when unethical conduct is detected.</p> <p>III. Detailed documentation of stakeholder collaboration and regulatory activities, including meeting records, workshop outcomes, and evidence of how collective input shapes interaction protocols.</p> <p>IV. Regular assessment reports showing framework effectiveness and adaptation, including records of regulatory body decisions, dispute</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
				resolutions, and updates made to address emerging technological and ethical considerations.
<p>G7.2 – Integration of Fairness Controls in AI Systems</p> <p>(Systems should possess robust fairness mechanisms integrated throughout their planning, decision-making, and operational processes to ensure respect for human life, rights, dignity, and universal values. Organizations should implement comprehensive frameworks that embed ethical principles and societal norms directly into AI system designs, preventing bias and discrimination while maintaining transparent and equitable operations)</p>	<p>a. Implement comprehensive ethical frameworks combining bias detection systems, fairness algorithms, and continuous training processes that ensure adherence to human rights and universal values while preventing discriminatory outcomes in decision-making.</p> <p>b. Establish multi-layered protection architecture incorporating safety protocols, transparency mechanisms, and monitoring systems that safeguard individual and community well-being while enabling clear oversight and timely human intervention.</p> <p>c. Foster inclusive development processes that involve diverse stakeholder groups in system design and evaluation, ensuring consideration of evolving societal values while promoting diversity in both development teams and training datasets.</p>	<p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Complete documentation of ethical frameworks and fairness mechanisms, including bias detection strategies, algorithmic fairness methodologies, and records demonstrating successful prevention of discriminatory outcomes.</p> <p>II. Comprehensive records of protection systems and oversight mechanisms, including safety protocols, transparency tools, monitoring logs, and evidence of effective human intervention capabilities.</p> <p>III. Detailed documentation of stakeholder engagement and diversity initiatives, including workshop records, survey results, and evidence of how diverse perspectives shape system design and improvement.</p> <p>IV. Regular assessment reports showing framework effectiveness and adaptation, including audit logs, compliance tests, and records of corrective actions taken to maintain alignment with ethical standards and societal values.</p>
<p>G7.3 – Balanced Global AI Partnership Framework</p> <p>(Systems should facilitate equitable distribution of AI capabilities and resources through balanced international partnerships. Organizations</p>	<p>a. Implement comprehensive international frameworks that enable equitable resource distribution and technology sharing while preventing dominance by powerful entities, with particular emphasis on including developing</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>I. Complete documentation of international frameworks and agreements, including technology sharing protocols, capacity building programs, and records demonstrating successful inclusion of developing nations in AI alliances.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>should establish frameworks that ensure fair technology sharing and knowledge exchange while actively preventing powerful entities from exploiting technological disparities or undermining global equilibrium through self-interested actions)</p>	<p>nations and marginalized groups in meaningful alliance participation.</p> <p>b. Establish transparent oversight mechanisms and governance structures that identify and prevent exploitative practices while ensuring diverse stakeholder participation in decision-making and accountability processes.</p> <p>c. Foster global education and collaborative research initiatives that enhance AI expertise worldwide, with particular focus on reducing technological disparities between developed and developing nations.</p>	<p>N</p> <p>I</p>	<p>D, I, O, M, U, R</p> <p>D, I, O, M, U, R</p>	<p>II. Comprehensive records of oversight activities and governance processes, including documentation of stakeholder participation, preventive measures against exploitation, and evidence of effective intervention against power imbalances.</p> <p>III. Detailed documentation of educational programs and research collaborations, including curricula, training materials, joint project outcomes, and impact assessments showing reduction in technological disparities.</p> <p>IV. Regular independent assessment reports evaluating framework effectiveness, including evidence of improved resource distribution, reduced disparities, and successful prevention of exploitative practices.</p>
<p>G7.4 – Control and Oversight of AI Autonomy</p> <p>(Systems should possess adaptable mechanisms that enable precise control over their degrees of autonomy while preventing improper interactions or exploitation. Organizations should implement comprehensive frameworks that integrate human oversight throughout decision-making processes while maintaining clear boundaries on autonomous operations)</p>	<p>a. Implement comprehensive control architecture combining adjustable autonomy levels, failsafe protocols, and human-in-the-loop systems that enable operators to modulate AI behavior based on performance metrics and risk assessments while ensuring rapid human intervention when needed.</p> <p>b. Establish rigorous monitoring frameworks incorporating continuous auditing, validation tools, and accountability logs that track both AI activities and human operator decisions while maintaining transparency in all autonomy-related adjustments.</p>	<p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Complete documentation of autonomy control frameworks, including technical specifications, operational parameters, and records demonstrating effective human modulation of AI behavior through whitelisting, blacklisting, and other control mechanisms.</p> <p>II. Comprehensive monitoring and audit records, including operator accountability logs, anomaly detection reports, and evidence of successful human intervention in high-risk scenarios or unexpected situations.</p> <p>III. Detailed documentation of ethical guidelines and compliance measures, including evidence of alignment with</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
	<p>c. Deploy embedded ethical and legal guidelines that ensure operations remain within authorized scopes while promoting compliance with societal norms and enabling clear understanding of AI decision-making processes.</p>	<p>N</p>	<p>D, I, O, M, R</p>	<p>societal norms and records showing consistent operation within authorized boundaries.</p> <p>IV. Regular assessment reports including case studies of failsafe protocol activation, human intervention outcomes, and evidence of effective oversight mechanisms in maintaining appropriate autonomy constraints.</p>
<p>G7.5 – Integration of AI Ethics Education (Systems should possess integrated mechanisms for promoting ethical awareness and understanding among developers and users through educational initiatives. Organizations should facilitate comprehensive AI ethics education that builds foundational competence in ethical implications, responsibilities, and impacts while fostering commitment to responsible AI development)</p>	<p>a. Establish collaborative frameworks between academic institutions, industry experts, and ethicists to develop standardized AI ethics curricula that combine technical knowledge with ethical principles, incorporating real-world case studies and practical insights into ethical decision-making.</p> <p>b. Foster interdisciplinary partnerships that enhance curriculum development through diverse perspectives while providing educators with ongoing professional development opportunities and updated resources to support effective ethics education.</p> <p>c. Extend ethics education beyond academia through community outreach and resource allocation that supports broad adoption of ethical practices in AI development and deployment.</p>	<p>N</p> <p>I</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<p>I. Complete documentation of educational partnerships and curriculum development, including meeting records, shared resources, and evidence of how diverse perspectives shape ethics education programs.</p> <p>II. Comprehensive records of interdisciplinary collaboration and educator support, including course materials, training programs, and evidence of continuous curriculum improvement based on emerging challenges.</p> <p>III. Detailed documentation of community outreach initiatives, including workshop agendas, participation metrics, and evidence of successful promotion of ethical practices beyond academic settings.</p> <p>IV. Regular assessment reports showing program effectiveness, including participant feedback, follow-up surveys, and evidence of increased ethical awareness and practice adoption among AI developers and users.</p>

Goal Title & Definition	AAI Safety Foundational Requirements (AAI-SFRs)	Normative/ Instructive	Stakeholder D, I, O, M, U, R	Required Evidence
<p>G7.6 – Integration of Human Ethics in AI Systems</p> <p>(Systems should possess deeply integrated ethical principles that enable them to autonomously uphold human rights and values throughout their decision-making processes. Organizations should implement comprehensive frameworks that ensure AI systems operate in harmony with human ethical norms while actively preventing the introduction of unintended biases during ethical training)</p>	<ul style="list-style-type: none"> a. Implement comprehensive ethical frameworks combining developer guidelines, universal human values, and bias detection mechanisms that ensure consistent ethical alignment while preventing unintended biases from emerging during training. b. Establish robust monitoring and explainability systems that enable continuous evaluation of ethical compliance while maintaining transparency in decision-making processes and facilitating effective human oversight. c. Foster sustained stakeholder engagement incorporating diverse perspectives, cultural sensitivity, and continuous learning mechanisms that enable adaptation to evolving societal norms and values. 	<p>N</p> <p>N</p> <p>N</p>	<p>D, I, O, M, R</p> <p>D, I, O, M, R</p> <p>D, I, O, M, R</p>	<ul style="list-style-type: none"> I. Complete documentation of ethical frameworks and developer guidelines, including training protocols, bias mitigation techniques, and records demonstrating successful alignment with human values and prevention of unintended biases. II. Comprehensive records of monitoring activities and oversight mechanisms, including audit reports, explainable AI methodologies, and evidence of effective detection and correction of ethical deviations. III. Detailed documentation of stakeholder consultation processes, including meeting records, feedback collection, and evidence of how diverse perspectives shape ethical guidelines and cultural sensitivity measures. IV. Regular assessment reports showing framework effectiveness and adaptation, including evidence of continuous learning processes and successful response to evolving societal norms.
<p>END</p>				