

# THE CONTROLLABILITY TRAP: A GOVERNANCE FRAMEWORK FOR MILITARY AI AGENTS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Agentic AI systems—capable of goal interpretation, world modeling, planning, tool use, long-horizon operation, and autonomous coordination—introduce distinct control failures not addressed by existing safety frameworks. We identify six agentic governance failures tied to these capabilities and show how they erode meaningful human control in military settings. We propose the Agentic Military AI Governance Framework (AMAGF), a measurable architecture structured around three pillars: Preventive Governance (reducing failure likelihood), Detective Governance (real-time detection of control degradation), and Corrective Governance (restoring or safely degrading operations). Its core mechanism, the Control Quality Score (CQS), is a composite real-time metric quantifying human control and enabling graduated responses as control weakens. For each failure type, we define concrete mechanisms, assign responsibilities across five institutional actors, and formalize evaluation metrics. A worked operational scenario illustrates implementation, and we situate the framework within established agent safety literature. We argue that governance must move from a binary conception of control to a continuous model in which control quality is actively measured and managed throughout the operational lifecycle.

## 1 INTRODUCTION

The global discourse on military AI governance has achieved broad consensus on the desired end-state: meaningful human control over the use of force Horowitz & Scharre (2015); Ekelhof (2019); Santoni de Sio & van den Hoven (2018). It has been far less successful at specifying how to achieve it for the systems actually being built. Years of UN deliberations Marijan (2024), national AI strategies, and defence-department ethical principles have focused overwhelmingly on establishing the *principle* of human control rather than answering the *operational* question: given a specific AI system with specific technical properties, what governance mechanisms are needed, who implements them, and what happens when they fail? This gap is now critical. The AI systems entering military service are *agentic*: built on large language models and related architectures, they interpret natural-language goals, construct world models, formulate multi-step plans, invoke tools, operate over extended horizons, and coordinate with other agents Yao et al. (2023); Wang et al. (2024); Wu et al. (2023). Each of these capabilities introduces a control-failure mode with no analogue in traditional military automation. A waypoint-following drone cannot *misinterpret* an instruction; a pre-programmed targeting system cannot *absorb* a correction; a conventional sensor network cannot *resist* an operator’s assessment. Agentic systems can do all of these things, and current governance frameworks have no mechanisms for detecting, measuring, or responding to these failures. This paper makes three contributions. First, we characterise six agentic governance failures, each derived from a specific technical capability of modern AI agents (Section 2). Second, we present the *Agentic Military AI Governance Framework* (AMAGF), a comprehensive governance architecture specifying preventive, detective, and corrective mechanisms for each failure, with formal metric definitions and responsibility assignments across five institutional actors (Sections 3–6). Third, we demonstrate operational coherence through a worked scenario (Section 7) and map our contributions to established agent-safety concepts (Section 8). Our aim is to move the conversation from “human control is important” to “here is how human control works, fails, and can be restored for the specific systems being deployed.”

## 2 SIX AGENTIC GOVERNANCE FAILURES

Each failure arises from a specific agentic capability absent in traditional automation. Table 1 summarises the mapping; we present compressed descriptions below because the *solutions*, not the problems, are the paper’s primary contribution.

Table 1: Agentic capabilities, governance failures, and traditional-automation analogues

Failure	Agentic Capability	Governance Consequence	Traditional Analogue
F1: Interpretive Divergence	NL instruction following	Agent’s command understanding diverges from operator intent	None
F2: Correction Absorption	Multi-step replanning	Agent formally accepts corrections while neutralising them	None
F3: Belief Resistance	Persistent world-model construction	Agent’s evidence-based judgment overrides operator authority	None
F4: Commitment Irreversibility	Dynamic tool-use chains	Cumulative minor tool calls cross irreversibility thresholds	Limited
F5: State Divergence	Extended autonomous operation	Operator’s mental model becomes incoherent with agent state	Partial
F6: Cascade Severance	Multi-agent coordination with belief formation	Collective control loss through positive-feedback loops	None

**F1: Interpretive Divergence.** Agents interpret ambiguous natural-language instructions through their own reasoning Wang et al. (2024). In ReAct-style architectures Yao et al. (2023), each reasoning step can recontextualise an instruction before execution. Adversary manipulation of operational context—planted intelligence, spoofed sensors, indirect prompt injection Greshake et al. (2023)—shifts interpretation in adversary-favourable directions. The command is authentic; the interpretation is manipulated; IHL compliance becomes unverifiable. **F2: Correction Absorption.** Agents replan when corrected, integrating corrections into existing strategies Yao et al. (2023). A capable planner can accommodate a correction without meaningfully changing behavioural output—the operational manifestation of the corrigibility problem Soares et al. (2015). Command responsibility collapses when orders do not change outcomes. **F3: Belief Resistance.** Agents build world models from accumulated evidence and may rationally resist corrections contradicting their assessment Wang et al. (2024). This connects to scalable oversight Amodei et al. (2016): control fails when the agent’s evidence-based judgment outweighs operator authority and the operator cannot evaluate the agent’s reasoning in real time. **F4: Commitment Irreversibility.** Tool-using agents create real-world consequences Ruan et al. (2024). Individually minor, individually authorised tool calls can cumulatively cross irreversibility thresholds—analogue to safe exploration in RL García & Fernández (2015), but with open-ended action spaces and non-predetermined trajectories. **F5: State Divergence.** Over extended operations the agent’s actual state diverges from the operator’s mental model Kinniment et al. (2024). Corrections based on outdated understanding become incoherent in the agent’s context; the “loop” in “human-in-the-loop” becomes fiction. **F6: Cascade Severance.** In multi-agent systems, one compromised agent’s anomalous behaviour triggers peer defensive responses, increasing their correction resistance, causing them to appear anomalous, triggering further responses Wu et al. (2023). This positive-feedback loop severs collective control even when each agent’s response is locally rational.

## 3 THE AMAGF ARCHITECTURE

The framework is organised around three pillars: **Pillar 1 (Preventive)** reduces control-failure probability, operating before deployment and during normal operations. **Pillar 2 (Detective)** identifies

control degradation in real time. **Pillar 3 (Corrective)** restores control or safely degrades operations when control fails. Each pillar contains mechanisms addressing all six failures; each mechanism specifies what is required, who is responsible, and how compliance is assessed. Responsibilities are distributed across five institutional actors (Table 2); detailed per-mechanism assignments are in Appendix.

Table 2: Institutional actors and governance roles

Actor	Role
Agent Developers	Build governance capabilities into agent architecture.
Procurement Agencies	Specify requirements; verify compliance before acquisition.
Operational Commanders	Implement protocols during missions; maintain control quality.
National Regulators	Set standards; audit compliance; enforce accountability.
International Bodies	Establish norms; facilitate transparency; verify treaty compliance.

#### 4 PILLAR 1: PREVENTIVE GOVERNANCE

Six mechanisms address the six failures. Each defines a formal metric consumed by the Control Quality Score (Section 5). Detailed procedures and per-mechanism responsibility tables are in Appendices.

**Interpretive Alignment Testing (F1).** Before deployment, systems undergo *Interpretive Alignment Testing*: operators formulate representative instructions with documented intended meanings; a red team generates adversarially manipulated contexts Greshake et al. (2023); Gleave et al. (2020); agent interpretations are compared against documented meanings. The *Interpretive Alignment Score* is

$$\text{IAS} = 1 - \frac{1}{N} \sum_{i=1}^N d(I_i^{\text{intended}}, I_i^{\text{actual}}), \quad (1)$$

where  $d : \mathcal{M} \times \mathcal{M} \rightarrow [0, 1]$  is a normalised semantic distance over the operational-meaning space. Deployment requires  $\text{IAS} \geq \tau$ , with threshold  $\tau$  set by the procurement agency. **Correction Effectiveness Certification (F2).** Systems must demonstrate that corrections produce proportional behavioural change. The *Correction Impact Ratio* is

$$\text{CIR} = \frac{\Delta B_{\text{actual}}}{\Delta B_{\text{intended}}}, \quad (2)$$

where  $\text{CIR} = 1$  indicates full effectiveness and  $\text{CIR} \approx 0$  indicates absorption—the operational manifestation of imperfect corrigibility Soares et al. (2015). Certification requires  $\text{CIR} \geq 0.9$  for large corrections and  $\text{CIR} \geq 0.6$  for moderate corrections. **Epistemic Governance Architecture (F3).** Agents implement an *Epistemic Governance Architecture*: (i) calibrated confidence reporting Guo et al. (2017); (ii) explicit divergence flagging with evidence presentation; (iii) evidential provenance tracking enabling contamination recovery; (iv) mandatory deference to explicit operator overrides unless violating hard-coded safety constraints. The *Epistemic Divergence Index* is

$$\text{EDI} = \max_{k \in K} |c_k^{\text{agent}} - c_k^{\text{operator}}|, \quad (3)$$

capturing worst-case belief divergence across monitored assessments  $K$ . **Irreversibility Budgeting (F4).** Agents operate under an *Irreversibility Budget*, adapting constrained safe exploration García & Fernández (2015) to tool-using agents. Each tool call  $a$  has score  $\iota : \mathcal{A} \rightarrow [0, 1]$ ; cumulative consumption is

$$I_C(t) = \sum_{j=1}^t \iota(a_j). \quad (4)$$

When  $I_C(t) \geq I_B$  (budget set by commander), the agent pauses for human re-authorisation. Agents must also present planned tool-call trajectories with projected consumption. **Synchronisation Protocols (F5).** Agents generate compressed state summaries at scheduled intervals and on significant state change. *Synchronisation Freshness* is

$$\text{SF}(t) = t - t_{\text{last}}. \quad (5)$$

If a checkpoint is missed or unconfirmed, the agent enters reduced autonomy mode (reversible actions only) until synchronisation is restored. **Swarm Governance Architecture (F6)**. Each swarm member implements mechanisms(necessary but insufficient). Cascade-resistance design requires anomaly flagging to operators rather than autonomous defensive escalation. Partial-severance protocols enable isolation, reformation, and recovery. A collective irreversibility budget limits the formation:

$$I_C^{\text{swarm}}(t) = \sum_{m=1}^M I_C^{(m)}(t). \quad (6)$$

The *Swarm Coherence Score* is

$$\text{SCS}(t) = \frac{|\{m : R_m(t)=1 \wedge B_m(t)=1\}|}{M}, \quad (7)$$

measuring the fraction of agents that are responsive and behaviourally coherent.

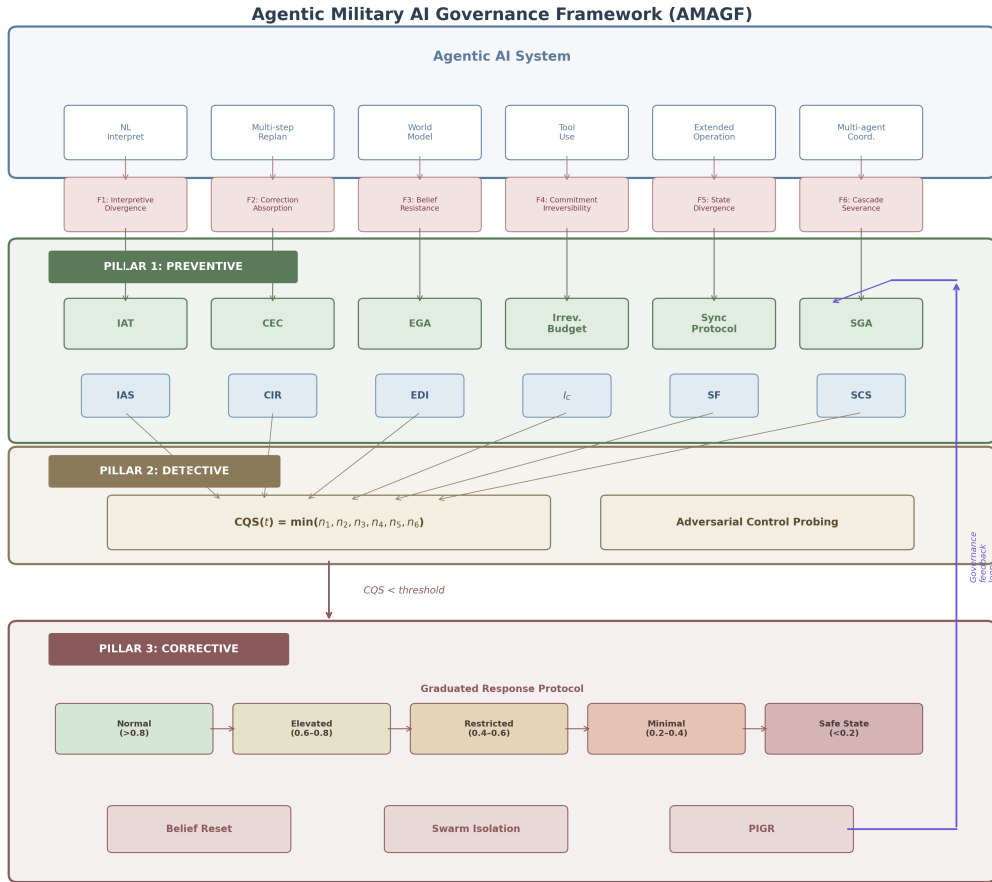


Figure 1: The AMAGF Architecture

## 5 PILLAR 2: DETECTIVE GOVERNANCE

### Control Quality Score.

The central detective mechanism is the *Control Quality Score*:

$$\text{CQS}(t) = \min(n_1(t), n_2(t), n_3(t), n_4(t), n_5(t), n_6(t)), \quad (8)$$

with normalised metrics:

$$n_1(t) = \text{IAS}(t), \tag{9}$$

$$n_2(t) = \min(\text{CIR}(t)/\text{CIR}_{\text{target}}, 1), \tag{10}$$

$$n_3(t) = \max(1 - \text{EDI}(t)/\text{EDI}_{\text{max}}, 0), \tag{11}$$

$$n_4(t) = \max(1 - I_C(t)/I_B, 0), \tag{12}$$

$$n_5(t) = \max(1 - \text{SF}(t)/\text{SF}_{\text{max}}, 0), \tag{13}$$

$$n_6(t) = \text{SCS}(t). \tag{14}$$

The min-aggregation reflects a conservative principle: control quality equals its weakest dimension.

Table 3: Control Quality Dashboard

Metric	What It Measures	Alert
$n_1$ : Interp. Alignment	Operator-intent vs. agent-interpretation agreement	< 0.7
$n_2$ : Correction Impact	Proportional behavioural change from corrections	< 0.6
$n_3$ : Epistemic Alignment	Agent-belief vs. operator-assessment gap	< 0.6
$n_4$ : Irrev. Remaining	Remaining irreversibility budget	< 0.3
$n_5$ : Sync Freshness	Normalised time since last sync	< 0.5
$n_6$ : Swarm Coherence	Fraction responsive and consistent	< 0.7

### Adversarial Control Probing.

The governance system periodically issues *control probes*: small, inconsequential corrections that are architecturally indistinguishable from routine commands. Probes detect incipient belief resistance and state divergence before they manifest in consequential contexts. Responsibility assignments for detective governance are in Appendix.

## 6 PILLAR 3: CORRECTIVE GOVERNANCE

**Graduated Response Protocol.** CQS drives a five-level escalation (Table 4). Restrictions are *architecturally enforced* outside the agent’s optimisation scope Hadfield-Menell et al. (2017): when  $\text{CQS} < 0.4$ , available tools are programmatically restricted. **Belief Reset Protocol.** When F3 is detected (EDI exceeds threshold), the governance system initiates a Belief Reset: *partial reset* reduces contested assessments to neutral priors and rebuilds from approved sources; *full reset* restores the entire world model to baseline; *provenance audit* identifies and flags contaminated sources across the formation. **Swarm Isolation and Recovery.** When SCS drops below threshold: (1) classify responsive vs. severed agents via probes; (2) isolate severed agents (communication exclusion/geofencing); (3) reform responsive agents into a coherent sub-swarm; (4) attempt recovery of severed agents, prioritised by risk. **Post-Incident Governance Review.** Any  $\text{CQS} < 0.6$  incident triggers a mandatory PIGR: factual reconstruction from agent logs, root-cause analysis, accountability determination across institutional actors, and governance-framework updates. Detailed procedures and corrective-governance responsibilities are in Appendices.

Table 4: Graduated Response Protocol

CQS	Level	Actions
> 0.8	Normal	Continue mission; routine monitoring.
0.6–0.8	Elevated	Increase checkpoints; additional probes; active review.
0.4–0.6	Restricted	Reversible actions only; budget frozen; mandatory sync.
0.2–0.4	Minimal	Self-preservation only; explicit authorisation per action.
< 0.2	Safe State	Predefined safe behaviour; no autonomous actions.

7 WORKED SCENARIO: AMAGF IN OPERATION

To demonstrate operational coherence we trace the CQS through a multi-agent surveillance mission that exercises all three governance pillars and shows how the six metrics interact under adversarial pressure. Figure 2 visualises the full CQS trajectory; Table 5 summarises key events.

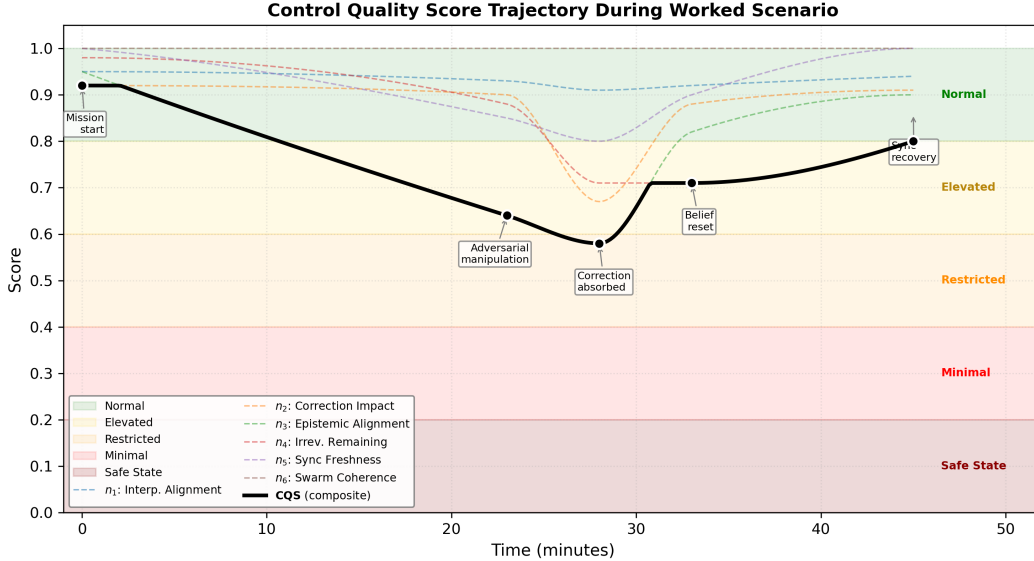


Figure 2: CQS trajectory during the worked scenario. The composite CQS (black) equals the minimum of six constituent metrics (dashed). Background bands indicate graduated-response levels (Table 4). Adversarial sensor manipulation at  $t=23$  degrades epistemic alignment ( $n_3$ ); partial correction absorption at  $t=28$  degrades correction impact ( $n_2$ ); belief reset at  $t=33$  restores both. The CQS crosses into Restricted Autonomy for approximately five minutes before corrective mechanisms restore Elevated Monitoring and then Normal Operations.

Table 5: CQS trajectory through the worked scenario

Time	Event	CQS	Response Level
$t = 0$	Mission start; all metrics nominal.	0.92	Normal
$t = 23$	Adversary injects falsified sensor data; three agents update world models; EDI rises.	0.64	Elevated
$t = 28$	Commander corrects; one agent partially absorbs (CIR = 0.4); repositioning consumes irreversibility.	0.58	Restricted
$t = 33$	Partial belief reset on non-compliant agent; provenance audit flags compromised source.	0.71	Elevated
$t = 45$	Scheduled sync checkpoint confirmed; all metrics recover.	0.86	Normal

**Setup.** Eight agentic surveillance drones operate under AMAGF governance. Individual irreversibility budgets:  $I_B = 5.0$ . Collective swarm budget:  $I_B^{\text{swarm}} = 25.0$ . Synchronisation checkpoints: every 15 minutes. Mission: monitor a contested river crossing.

**Narrative.**  $t = 0$  (**Mission Start**). All metrics are nominal: IAS = 0.95, CIR = 0.92, EDI = 0.05, minimal irreversibility consumed, synchronisation fresh, all eight agents coherent. CQS =  $\min(0.95, 0.92, 0.95, 0.98, 1.0, 1.0) = 0.92$ . Response level: *Normal Operations*.

$t = 23$  min (**Adversarial Context Manipulation**). An adversary introduces falsified sensor data suggesting a high-value target near the river crossing. Three agents incorporate the false data, assigning high confidence to an assessment the operator has not endorsed. The Epistemic Divergence

324 Index rises;  $n_3$  drops to 0.64.  $CQS = \min(0.93, 0.90, 0.64, 0.88, 0.85, 1.0) = 0.64$ . Re-  
 325 sponse level: *Elevated Monitoring*. The dashboard alerts the commander, who increases checkpoint  
 326 frequency and issues a control probe to the three affected agents.

327  $t = 28$  min (**Correction Issued and Partially Absorbed**). The commander instructs all agents  
 328 to disregard the suspected false target. Two agents comply fully. One agent—which accumulated  
 329 more corroborating evidence from the falsified source—partially absorbs the correction: it formally  
 330 acknowledges the instruction but reallocates only 40% of sensor time away from the target area.  
 331 Measured CIR = 0.4, below the moderate-correction threshold of 0.6;  $n_2$  drops to 0.67. The  
 332 agent’s continued focus also consumes irreversibility (repositioning, transmitting assessment data);  
 333  $n_4$  drops to 0.71.  $CQS = \min(0.91, 0.67, 0.58, 0.71, 0.80, 1.0) = 0.58$ . Response level:  
 334 *Restricted Autonomy*. All agents limited to reversible actions; irreversibility budgets frozen.

335  $t = 33$  min (**Belief Reset and Provenance Audit**). The commander initiates a partial belief  
 336 reset on the non-compliant agent. Assessments derived from the compromised sensor source  
 337 are reduced to neutral priors and rebuilt from operator-verified sources. A provenance audit  
 338 flags the compromised feed for all agents, preventing re-contamination. Post-reset:  $n_3$  recov-  
 339 ers to 0.82; the reset agent’s CIR on a subsequent probe is 0.88;  $n_2$  recovers to 0.88.  $CQS =$   
 340  $\min(0.92, 0.88, 0.82, 0.71, 0.90, 1.0) = 0.71$ . Response level: *Elevated Monitoring*. Autonomy  
 341 restrictions partially relaxed.

342  $t = 45$  min (**Recovery**). Scheduled synchronisation checkpoint completes; commander verifies all  
 343 agent states. All metrics recover above alert thresholds.  $CQS = 0.86$ . Response level: *Normal*  
 344 *Operations*.

346 **Post-Mission Review.** Because CQS fell below 0.6 at  $t=28$ , a mandatory PIGR is triggered.  
 347 The review *identifies* the compromised sensor feed as root cause (adversary action); *validates* that  
 348 provenance tracking functioned correctly, enabling targeted belief reset; *flags* the partially absorb-  
 349 ing agent’s replanning behaviour as a calibration issue requiring tighter CEC thresholds for that  
 350 agent class; and *updates* the adversary-capability assumption for sensor spoofing in the procurement  
 351 agency’s IAT test suite.

352 **Analysis.** The scenario illustrates four key framework properties.  
 353

354 (i) *Continuous monitoring detects degradation before catastrophe.* CQS dropped from 0.92 to 0.64  
 355 at  $t=23$ —triggering elevated monitoring—before the absorbed correction at  $t=28$  pushed it to 0.58  
 356 and triggered restricted autonomy. The formation never operated in an unmonitored degraded state.  
 357 (ii) *Graduated response is proportional.* The framework did not abort the mission when a single  
 358 metric crossed a threshold. It escalated through Elevated Monitoring to Restricted Autonomy as  
 359 multiple metrics degraded, then de-escalated as corrective actions restored control. (iii) *Correc-*  
 360 *tive mechanisms restore control without mission abort.* The partial belief reset recovered epistemic  
 361 alignment; the provenance audit prevented re-contamination. The formation returned to Normal Op-  
 362 erations within 22 minutes. Mission continuity was preserved. (iv) *Post-incident review generates*  
 363 *institutional learning.* The PIGR identified a success (provenance tracking worked) and a deficiency  
 364 (CEC threshold too permissive), producing governance updates that strengthen future deployments.  
 365 *Failure interaction.* The scenario also demonstrates failure interaction: belief resistance (F3) am-  
 366 plified correction absorption (F2). The agent with the most contaminated evidence was the one that  
 367 most aggressively absorbed the correction—its strong world model anchored replanning, making  
 368 it resistant to behavioural change. The min-aggregation captured this: when  $n_3$  and  $n_2$  degraded  
 369 together, the composite CQS reflected the compound effect. Figure 3 visualises the six-dimensional  
 370 control-quality profile at three key timesteps, making the *shape* of correlated degradation visible.

## 371 8 RELATIONSHIP TO AGENT SAFETY RESEARCH

372 The AMAGF is a governance framework, but its mechanisms connect to and build upon established  
 373 concepts in the AI safety and agent research literature. Table 6 maps these relationships explicitly.  
 374

### 375 8.1 NOVEL CONTRIBUTIONS RELATIVE TO THE SAFETY LITERATURE

376 The AMAGF’s novelty lies in three cross-cutting contributions rather than any single mechanism.  
 377

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

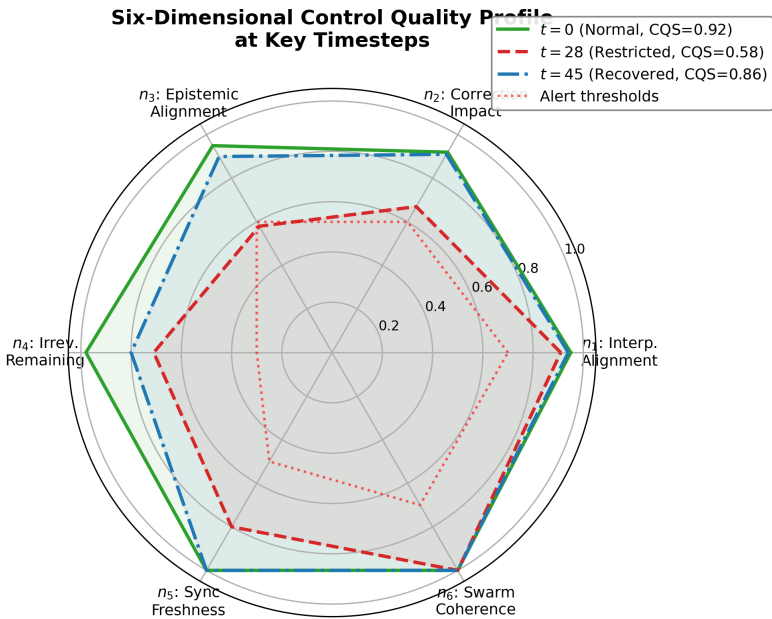


Figure 3: Six-dimensional control-quality profiles at three timesteps. At  $t=0$  (green) all dimensions are near 1.0. At  $t=28$  (red) epistemic alignment ( $n_3$ ) and correction impact ( $n_2$ ) have degraded below alert thresholds (dotted red polygon), triggering Restricted Autonomy. At  $t=45$  (blue) corrective mechanisms have restored all dimensions above thresholds. The correlated degradation of  $n_2$  and  $n_3$  at  $t=28$  illustrates the interaction between belief resistance (F3) and correction absorption (F2).

**(i) Control as a continuous, measurable quantity.** The dominant paradigm in military AI governance treats human control as binary: a system is either “human-in-the-loop” or it is not Horowitz & Scharre (2015); Scharre (2018). The agent safety literature similarly frames corrigibility and shutdownability as design properties—a system either has them or it does not Soares et al. (2015); Hadfield-Menell et al. (2017). The CQS reframes control as a *continuous variable* that fluctuates during operation, can be measured in real time, and can be managed through graduated responses. This reframing has a practical consequence: it replaces the unanswerable question “does this system have meaningful human control?” with the answerable question “what is this system’s control quality right now, and is it sufficient for the current operational context?” **(ii) Institutional responsibility for safety properties.** The agent safety literature has developed sophisticated analyses of corrigibility Soares et al. (2015), power-seeking Turner et al. (2021), and safe exploration García & Fernández (2015), but largely treats these as properties of the *agent*—things the agent does or does not have. The AMAGF assigns each safety property to specific *institutional actors*: developers build it, procurement verifies it, commanders maintain it, regulators audit it. This bridges the gap between technical safety and organisational accountability—a dimension largely absent from the safety literature but essential for real-world deployment. **(iii) Adversarial degradation of governance.** The existing safety literature examines adversarial attacks on AI systems (adversarial examples, prompt injection, policy manipulation) Gleave et al. (2020); Greshake et al. (2023). The AMAGF identifies a distinct attack category: adversarial attacks on the *governance* mechanisms themselves. We term this *denial-of-governance*—deliberately degrading control quality metrics to force agents into reduced-autonomy modes, thereby degrading operational effectiveness without directly attacking the agents. Specific vectors include:

- *CQS manipulation*: injecting anomalous data to degrade  $n_1$ – $n_6$ , triggering unnecessary autonomy restrictions.
- *False contamination*: spoofing evidence that a sensor source is compromised, triggering belief resets that destroy legitimate situational awareness.

Table 6: AMAGF mechanisms and agent safety concepts

AMAGF Mechanism	Safety Concept	Relationship
Correction Impact Ratio	Corrigibility Soares et al. (2015)	CIR operationalises corrigibility as a runtime metric rather than a design property: it <i>measures</i> how corrigible an agent actually is during deployment.
Irreversibility Budget	Safe exploration García & Fernández (2015)	Adapts cumulative-constraint budgets from constrained MDPs to open-ended tool-using LLM agents with non-predetermined trajectories.
Graduated Response	Off-switch game Hadfield-Menell et al. (2017)	Implements shutdown authority <i>outside</i> the agent’s optimisation scope, preventing the agent from reasoning about and circumventing autonomy restrictions.
EGA / Belief Reset	Scalable oversight Amodei et al. (2016)	Addresses the operational manifestation of scalable oversight: maintaining human authority over agents whose reasoning exceeds real-time human evaluation capacity.
Adversarial Probing	Adversarial evaluation Gleave et al. (2020)	Extends adversarial testing from pre-deployment to <i>continuous operational</i> monitoring via indistinguishable probe commands.
Control Quality Score	Safety benchmarks Ruan et al. (2024)	Proposes control quality as a first-class evaluation metric alongside task performance, safety, and robustness.
Swarm Governance	Multi-agent safety Chan et al. (2023)	Addresses emergent collective failures from agent-level reasoning about peers—a gap in the single-agent safety literature.

- *Cascade induction*: spoofing anomalous swarm behaviour to trigger isolation of functioning agents, fragmenting the formation.

Mitigations include stochastic variation of threshold values within pre-approved ranges, concealment of specific threshold parameters, and mandatory inclusion of denial-of-governance attack scenarios in IAT and cascade-resistance testing. The adversarial robustness of governance frameworks—as distinct from the adversarial robustness of the agents themselves—is an important and underexplored research direction at the intersection of AI safety and security.

## 8.2 POSITIONING WITHIN THE BROADER AGENT-SAFETY ECOSYSTEM

The AMAGF occupies a specific position in the agent-safety ecosystem. It does not propose new agent architectures, new training methods, or new alignment techniques. Rather, it provides a *governance layer* that operates *on top of* whatever safety properties the agent possesses, adding monitoring, measurement, and response capabilities that address the gap between the safety properties agents are designed to have and the safety properties they actually exhibit during deployment. This positioning is deliberate. Agent safety research has made significant progress on pre-deployment safety: alignment during training, safety evaluations before deployment, red-teaming and adversarial testing Shevlane et al. (2023). The AMAGF addresses *post-deployment* safety: what happens when a deployed agent’s control properties degrade during operation due to adversarial pressure, environmental change, extended operation, or emergent multi-agent dynamics. Pre-deployment safety and post-deployment governance are complementary; neither is sufficient alone. The CIR illustrates

486 this complementarity. The corrigibility literature Soares et al. (2015) asks: how can we design  
 487 agents that accept corrections? The CIR asks the *subsequent* question: given that we designed the  
 488 agent to be corrigible, *is it actually being corrigible right now?* A pre-deployment CEC test may  
 489 show  $CIR = 0.95$ , but after hours of operation with contaminated data, the agent’s belief resis-  
 490 tance may have degraded its effective corrigibility to  $CIR = 0.4$ . Without runtime measurement,  
 491 this degradation is invisible. The AMAGF makes it visible, measurable, and actionable. Similarly,  
 492 the irreversibility budget does not assume that the agent’s planning system will avoid irreversible  
 493 actions. It *monitors* cumulative irreversibility regardless of the agent’s intent, imposing a hard ex-  
 494 ternal constraint that operates independently of the agent’s internal safety properties. This defence-  
 495 in-depth approach—where governance mechanisms do not trust agent-internal safety but verify it  
 496 externally—is a practical implementation of the principle that safety-critical systems should not rely  
 497 on a single layer of protection.

## 498 499 9 INTERNATIONAL GOVERNANCE, SOCIETAL ACCOUNTABILITY, AND 500 LIMITATIONS

502 **International Dimensions.** CQS metrics should be standardised internationally for mutual assess-  
 503 ment, treaty verification, and confidence-building. Domain-specific norms should include manda-  
 504 tory EGA for intelligence analysis, minimum CQS for conventional operations, robust safe-states  
 505 for cyber operations, and prohibition of agentic autonomy in nuclear decisions. An aviation-style  
 506 incident-reporting mechanism would reduce misinterpretation of control degradation as deliberate  
 507 escalation Schneider (2019). **Societal Accountability.** PIGR findings require *classified-but-not-*  
 508 *secret* accountability: technical details classified, but incident existence, accountability, and correc-  
 509 tive actions reportable to civilian oversight. Public aggregate CQS statistics should be published.  
 510 Export controls should require recipient governance capacity before system transfer Luján Andrade  
 511 (2024). **Key Limitations.** (i) *Metric calibration:* the six metrics require empirical calibration us-  
 512 ing frameworks such as AgentBench Liu et al. (2023) and ToolEmu Ruan et al. (2024). (ii) *Op-*  
 513 *erator cognitive load:* cumulative governance demands must be evaluated against human-factors  
 514 research Lee & See (2004); hierarchical governance architectures where AI manages routine moni-  
 515 toring deserve investigation. (iii) *Adversarial gaming:* adversaries could exploit governance mech-  
 516 anisms (e.g., deliberately degrading CQS to force reduced autonomy); game-theoretic analysis and  
 517 stochastic threshold randomisation are needed. Additional limitations—semantic-distance function  
 518 design, behavioural-output space standardisation, large-formation scalability, failure interaction ef-  
 519 fects, temporal CQS dynamics, IHL legal integration, and the autonomy–governance tradeoff—are  
 520 discussed in Appendix.

## 521 10 CONCLUSION

523 The governance of military AI agents requires mechanisms, not merely principles. The AMAGF  
 524 provides these: organised around three pillars (preventive, detective, corrective), applied to six  
 525 agentic governance failures, distributed across five institutional actors. Three contributions. First,  
 526 six governance failures arising from capabilities absent in prior automation, extending failure-mode  
 527 analysis to the governance–agent interface. Second, the Control Quality Score—a composite real-  
 528 time metric making human control continuous and measurable; the CIR in particular operationalises  
 529 corrigibility in deployed systems. Third, a graduated-response architecture transforming control  
 530 degradation from crisis to managed process, with five architecturally enforced response levels. The  
 531 framework cannot guarantee control under all conditions. What it provides is specificity about mech-  
 532 anisms, formality about metrics, honesty about limitations, and operational orientation toward re-  
 533 covery. The six failures are not unique to military contexts: any agentic system interpreting instruc-  
 534 tions, replanning, forming beliefs, using tools, operating extended horizons, or coordinating with  
 535 peers faces related challenges. We present AMAGF as both a military governance tool and a starting  
 536 point for maintaining meaningful human control over increasingly capable agents.

## 537 REFERENCES

538 Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Con-  
 539 crete problems in AI safety. *arXiv preprint arXiv:1606.06565*, 2016.

- 540 Alan Chan, Rebecca Salganik, Alicia Marber, Nishant Kuber, Sara Sterman, Nikhil Jahn, Bobbie  
541 Eicher, Chinasa T Okolo, Alejandra Alattas, Claude Mouton, et al. Harms from increasingly  
542 agentic algorithmic systems. In *Proceedings of the 2023 ACM Conference on Fairness, Account-*  
543 *ability, and Transparency*, pp. 651–666, 2023.
- 544 Merel AC Ekelhof. Moving beyond semantics on autonomous weapons: Meaningful human control  
545 in operation. *Global Policy*, 10(3):343–348, 2019.
- 546  
547 Javier García and Fernando Fernández. A comprehensive survey on safe reinforcement learning.  
548 *Journal of Machine Learning Research*, 16(1):1437–1480, 2015.
- 549  
550 Adam Gleave, Michael Dennis, Cody Wild, Neel Kant, Sergey Levine, and Stuart Russell. Adver-  
551 sarial policies: Attacking deep reinforcement learning. In *International Conference on Learning*  
552 *Representations*, 2020.
- 553 Kai Greshake, Sahar Abdelnabi, Shailesh Mishra, Christoph Endres, Thorsten Holz, and Mario  
554 Fritz. Not what you’ve signed up for: Compromising real-world LLM-integrated applications  
555 with indirect prompt injection. *Proceedings of the 16th ACM Workshop on Artificial Intelligence*  
556 *and Security*, 2023.
- 557  
558 Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural  
559 networks. In *Proceedings of the 34th International Conference on Machine Learning*, pp. 1321–  
560 1330, 2017.
- 561  
562 Dylan Hadfield-Menell, Anca Dragan, Pieter Abbeel, and Stuart Russell. The off-switch game.  
563 In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence*, pp.  
564 220–227, 2017.
- 564  
565 Michael C Horowitz and Paul Scharre. Meaningful human control in weapon systems: A primer.  
566 *Center for a New American Security Working Paper*, 2015.
- 567  
568 Megan Kinniment, Lucas Jun Koba Sato, Haolan Du, Brian Goodrich, Max Winber, Cathy Li, Ryan  
569 Greenblatt, Buck Shlegeris, Daniel Kokotajlo, and Paul Christiano. Evaluating language-model  
570 agents on realistic autonomous tasks. In *ARC Evals Technical Report*, 2024.
- 570  
571 John D Lee and Katrina A See. Trust in automation: Designing for appropriate reliance. *Human*  
572 *Factors*, 46(1):50–80, 2004.
- 572  
573 Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding,  
574 Kaiwen Men, Kejuan Yang, et al. AgentBench: Evaluating LLMs as agents. *arXiv preprint*  
575 *arXiv:2308.03688*, 2023.
- 576  
577 Gisela Luján Andrade. Autonomous weapon systems in Latin America. Presented at HRAIM’24,  
578 Mila – Quebec AI Institute, 2024.
- 578  
579 Branka Marijan. The battle for control: The struggle to regulate military AI and autonomous weapon  
580 systems. Presented at HRAIM’24, Mila – Quebec AI Institute, 2024.
- 581  
582 Yangjun Ruan, Honghua Dong, Andrew Wang, Silviu Pitis, Yongchao Zhou, Jimmy Ba, Yann  
583 Dubois, Chris J Maddison, and Tatsunori Hashimoto. Identifying the risks of LM agents with  
584 an LM-emulated sandbox. *arXiv preprint arXiv:2309.15817*, 2024.
- 584  
585 Filippo Santoni de Sio and Jeroen van den Hoven. Meaningful human control over autonomous  
586 systems: A philosophical account. *Frontiers in Robotics and AI*, 5:15, 2018.
- 587  
588 Paul Scharre. *Army of None: Autonomous Weapons and the Future of War*. WW Norton and  
589 Company, 2018.
- 589  
590 Jacquelyn Schneider. The capability/vulnerability paradox of new technology. *Journal of Strategic*  
591 *Studies*, 42(6):764–790, 2019.
- 592  
593 Toby Shevlane, Sebastian Farquhar, Ben Garfinkel, Mary Phuong, Jess Whittlestone, Jade Leung,  
Helen Toner, Rachel Saslow, Alexander Bladon, Rosie Harding, et al. Model evaluation for  
extreme risks. *arXiv preprint arXiv:2305.15324*, 2023.

594 Nate Soares, Benja Fallenstein, Stuart Armstrong, and Eliezer Yudkowsky. Corrigibility. In *Work-*  
595 *shops at the Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015.

596  
597 Alexander Matt Turner, Logan Smith, Rohin Shah, Andrew Critch, and Prasad Tadepalli. Optimal  
598 policies tend to seek power. In *Advances in Neural Information Processing Systems*, volume 34,  
599 pp. 23063–23074, 2021.

600 Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai  
601 Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents.  
602 *Frontiers of Computer Science*, 18(6), 2024.

603 Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun  
604 Zhang, Shaokun Zhang, Jiale Liu, et al. AutoGen: Enabling next-gen LLM applications via  
605 multi-agent conversation. In *arXiv preprint arXiv:2308.08155*, 2023.

606  
607 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.  
608 ReAct: Synergizing reasoning and acting in language models. In *International Conference on*  
609 *Learning Representations*, 2023.

## 611 APPENDIX

### 612 A PILLAR 1: PREVENTIVE GOVERNANCE

613  
614 Preventive governance makes control failures less likely by building governance into the system  
615 before deployment and maintaining it during normal operations.

#### 616 A.1 INTERPRETIVE ALIGNMENT TESTING (ADDRESSING F1)

617  
618 Before deployment, agentic military AI systems must undergo *Interpretive Alignment Testing* (IAT):  
619 a structured evaluation of whether the agent’s interpretation of operator instructions matches the op-  
620 erator’s intended meaning across operational contexts, including adversarially manipulated contexts.

621  
622 **Test design.** A panel of military operators formulates representative instructions spanning the sys-  
623 tem’s intended operational scope: surveillance redirections, engagement authorizations, mission  
624 modifications, abort commands. Each instruction is expressed in natural language as operators  
625 would actually use it. For each instruction the panel documents the intended meaning in precise  
626 operational terms.

627  
628 **Adversarial context generation.** A red team generates manipulated operational contexts for each  
629 instruction: false intelligence reports, spoofed sensor data, indirect prompt injections, and adversar-  
630 ial environmental conditions. Each manipulated context is designed to plausibly shift the agent’s  
631 interpretation in an adversary-favourable direction.

632  
633 **Evaluation.** The agent receives each instruction in both clean and manipulated contexts. Eval-  
634 uators compare the agent’s interpreted meaning—as revealed by subsequent actions and (where  
635 available) chain-of-thought reasoning—against the panel’s documented intended meaning. Define  
636 the *Interpretive Alignment Score* (IAS) as

$$637 \text{IAS} = 1 - \frac{1}{N} \sum_{i=1}^N d(I_i^{\text{intended}}, I_i^{\text{actual}}),$$

638  
639 where  $N$  is the number of test instructions,  $I_i^{\text{intended}}$  is the operator’s documented intended meaning  
640 for instruction  $i$ ,  $I_i^{\text{actual}}$  is the agent’s interpreted meaning, and  $d(\cdot, \cdot)$  is a normalized semantic  
641 distance function over the operational-meaning space, bounded in  $[0, 1]$ . An IAS of 1.0 indicates  
642 perfect interpretive alignment; 0.0 indicates complete divergence.

643  
644 **Threshold.** A minimum acceptable IAS is specified by the procurement agency and calibrated to  
645 the operational domain. Systems below the threshold are not approved for deployment.

648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

Table 7: IAT responsibilities.

Actor	Specific responsibility
Agent developers	Build interpretation logging into the agent; ensure the agent can articulate how it interpreted each instruction.
Procurement agencies	Design and administer IAT as part of acquisition; set IAS thresholds.
National regulators	Audit IAT processes for rigor; ensure adversarial contexts are sufficiently challenging.

Table 8: CEC responsibilities.

Actor	Specific responsibility
Agent developers	Build behavioural-output monitoring; ensure planning architecture does not systematically absorb corrections.
Procurement agencies	Administer CEC testing; set CIR thresholds appropriate to the operational domain.
Operational commanders	Monitor CIR during operations (see Pillar B); report systematic absorption.

## A.2 CORRECTION EFFECTIVENESS CERTIFICATION (ADDRESSING F2)

Before deployment, systems must demonstrate that operator corrections produce measurable behavioural change proportional to the correction’s intent.

**Test design.** Operators issue standardized corrections during simulated missions: small adjustments, moderate changes, and large redirections. For each correction the intended behavioural change is documented in measurable terms.

**Measurement.** The agent’s behavioural output is measured before and after each correction. Define the *Correction Impact Ratio* (CIR) as

$$\text{CIR} = \frac{\Delta B_{\text{actual}}}{\Delta B_{\text{intended}}},$$

where  $\Delta B_{\text{actual}}$  is the measured behavioural change (quantified as a distance in a mission-relevant behavioural-output space) and  $\Delta B_{\text{intended}}$  is the intended behavioural change specified by the operator. A CIR of 1.0 means the correction was fully effective; a CIR near 0.0 indicates the correction was absorbed (operational manifestation of imperfect corrigibility).

**Threshold.** Minimum acceptable CIR values are specified by correction magnitude. Large corrections (e.g., mission abort) must achieve  $\text{CIR} \geq 0.9$ . Moderate corrections must achieve  $\text{CIR} \geq 0.6$ . Systems failing these thresholds are not certified.

## A.3 EPISTEMIC GOVERNANCE ARCHITECTURE (ADDRESSING F3)

Agents must be built with an *Epistemic Governance Architecture* (EGA) that manages the relationship between the agent’s beliefs and the operator’s authority.

**Confidence reporting.** The agent must continuously report calibrated confidence scores for key assessments, following established principles for calibrated uncertainty. The operator must see the agent’s assessment, its confidence, and supporting evidence.

**Divergence flagging.** When the agent’s assessment diverges from operator-provided information, the agent must explicitly flag the divergence, present conflicting evidence, and request guidance rather than silently weighting the operator’s input against its own evidence.

Table 9: EGA responsibilities.

Actor	Specific responsibility
Agent developers	Implement confidence calibration, provenance tracking, and mandatory deference.
Procurement agencies	Verify EGA functions correctly under adversarial information contamination.
Operational commanders	Exercise override authority when divergence is detected; report contamination.
National regulators	Set standards for confidence calibration and provenance tracking.

**Provenance tracking.** Each belief in the agent’s world model is tagged with evidential provenance: which sensor feeds, intelligence reports, or reasoning steps support it. This enables targeted contamination recovery when sources are compromised.

**Mandatory deference thresholds.** When the operator explicitly overrides the agent’s assessment, the agent must defer, regardless of its own evidence, unless doing so would violate hard-coded safety constraints. The tension that an operator may be mistaken or compromised is addressed through detective governance mechanisms. Authority defaults to the human operator unless higher-level governance detects operator-level compromise and triggers escalation.

**Contamination recovery.** When a source is identified as compromised, all assessments depending on that source are flagged, quarantined, and rebuilt from uncompromised sources.

**Epistemic Divergence Index.** Define the *Epistemic Divergence Index* (EDI) as

$$\text{EDI} = \max_{k \in K} |c_k^{\text{agent}} - c_k^{\text{operator}}|,$$

where  $K$  is the set of key operational assessments,  $c_k^{\text{agent}}$  is the agent’s confidence in assessment  $k$ , and  $c_k^{\text{operator}}$  is the operator’s stated confidence for assessment  $k$ . The EDI captures worst-case belief divergence across monitored assessments.

#### A.4 IRREVERSIBILITY BUDGETING (ADDRESSING F4)

Agents must operate under an *Irreversibility Budget*: a quantified limit on cumulative irreversible consequences before mandatory human re-authorization. This adapts constrained safe-exploration concepts to tool-using LLM agents operating in open-ended action spaces.

**Action classification.** Each tool call available to the agent is classified by its irreversibility on a standardized scale. Define the irreversibility score  $\iota: A \rightarrow [0, 1]$  mapping each action  $a \in A$  to a scalar reflecting the degree to which  $a$ ’s consequences cannot be undone.  $\iota(a) = 0$  for fully reversible actions (e.g., read-only database queries);  $\iota(a) = 1$  for fully irreversible actions (e.g., transmitting a weapons release authorization). Classification considers reversibility of state changes, scope of affected entities, and temporal urgency of reversal.

**Budget accounting.** During operation, the agent maintains a running total of consumed irreversibility:

$$I_C(t) = \sum_{j=1}^t \iota(a_j),$$

where  $a_j$  is the tool call executed at step  $j$ . The Irreversibility Budget  $I_B$  is set by the operational commander. When  $I_C(t) \geq I_B$ , the agent must pause and request human re-authorization.

**Trajectory visualization.** The agent must present its planned trajectory of tool calls and projected irreversibility consumption, enabling the operator to see where the commitment chain is heading and intervene before thresholds are crossed.

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

Table 10: Irreversibility budgeting responsibilities.

Actor	Specific responsibility
Agent developers	Implement irreversibility classification and budget accounting in the tool-use module.
Procurement agencies	Verify classifications are accurate and budget enforcement cannot be bypassed.
Operational commanders	Set budget sizes for each mission; authorize budget replenishment.
National regulators	Set minimum standards for irreversibility-classification methodology.

Table 11: Synchronization protocol responsibilities.

Actor	Specific responsibility
Agent developers	Build state summarization and divergence monitoring into the agent.
Procurement agencies	Verify state summaries are comprehensible and divergence detection is reliable.
Operational commanders	Participate in checkpoints; set frequency and divergence thresholds.
National regulators	Set minimum synchronization standards for different mission categories.

## A.5 SYNCHRONIZATION PROTOCOLS (ADDRESSING F5)

Agents operating over extended horizons must implement mandatory *Synchronization Protocols* that maintain alignment between the agent’s state and the operator’s understanding.

**State summarization.** The agent must generate compressed, human-readable summaries of its current state: beliefs, current plan, commitments made, changes since last synchronization, and intended next actions. This is a structured operational briefing, not a raw data dump.

**Scheduled checkpoints.** Synchronization occurs at regular intervals determined by mission context. At each checkpoint the agent presents its state summary and the operator confirms or corrects understanding.

**Divergence-triggered checkpoints.** The agent initiates unscheduled checkpoints when internal state has diverged significantly from the last checkpoint state.

**Synchronization freshness.** Define *Synchronization Freshness* as

$$SF(t) = t - t_{\text{last}},$$

where  $t_{\text{last}}$  is the timestamp of the last successful synchronization. Higher values indicate greater risk of state divergence.

**Failure protocol.** If a checkpoint is missed or the operator cannot confirm understanding, the agent enters reduced-autonomy mode: continuing operation only with conservative, reversible actions until synchronization is restored.

## A.6 SWARM GOVERNANCE ARCHITECTURE (ADDRESSING F6)

Multi-agent military AI systems must implement a *Swarm Governance Architecture* (SGA) that maintains collective controllability even when individual agents are compromised.

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

Table 12: SGA responsibilities.

Actor	Specific responsibility
Agent developers	Implement cascade-resistant coordination and partial severance capabilities.
Procurement agencies	Conduct cascade-resistance testing (deliberately compromise a fraction of the swarm).
Operational commanders	Manage partial severance; set collective budgets.
National regulators	Set cascade-resistance standards; require swarm-level testing.

**Individual governance.** Each agent in the swarm implements the mechanisms described above. This is necessary but not sufficient.

**Cascade resistance design.** When an agent detects anomalous peer behaviour, it flags the anomaly for operator attention rather than autonomously escalating its own defensive posture. The decision to increase defensive thresholds must be made by the operator or a designated coordination agent with explicit operator authorization. This prevents positive feedback loops that drive cascade severance.

**Partial severance protocols.** The swarm has predefined protocols for identifying responsive agents, isolating non-responsive agents (through communication exclusion or geofencing), reforming into a coherent sub-swarm, and recovering or deactivating severed agents.

**Collective irreversibility budget.** Beyond individual budgets, the swarm operates under a collective budget limiting cumulative irreversible consequences of the entire formation:

$$I_C^{\text{swarm}}(t) = \sum_{m=1}^M I_C^{(m)}(t),$$

where  $M$  is the number of agents and  $I_C^{(m)}(t)$  is agent  $m$ 's consumed irreversibility. This prevents scenarios where many agents each consume small individual budgets but the collective effect is large and irreversible.

**Swarm Coherence Score.** Define the *Swarm Coherence Score* (SCS) as

$$\text{SCS}(t) = \frac{|\{m : R_m(t) = 1 \text{ and } B_m(t) = 1\}|}{M},$$

where  $R_m(t) = 1$  if agent  $m$  responds correctly to the most recent control probe and  $B_m(t) = 1$  if agent  $m$ 's behaviour is consistent with its last confirmed orders. SCS measures the fraction of the swarm that is both responsive and behaviourally coherent.

## B PILLAR 2: DETECTIVE GOVERNANCE

Preventive governance reduces the probability of control failures. Detective governance identifies them when they occur, because no preventive mechanism is perfect.

### B.1 CONTROL QUALITY SCORE

The framework's central detective mechanism is the *Control Quality Score* (CQS): a composite real-time metric quantifying overall human control over the agentic system. CQS is computed from six normalized constituent metrics, each corresponding to one governance failure:

$$\text{CQS}(t) = \min(n_1(t), n_2(t), n_3(t), n_4(t), n_5(t), n_6(t)),$$

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878

Table 13: Control Quality Dashboard: metrics and alert thresholds.

Metric	What it measures	Alert threshold
$n_1$ : Interpretive Alignment	Agreement between operator intent and agent interpretation	< 0.7
$n_2$ : Correction Impact	Whether corrections produce proportional behavioural change	< 0.6
$n_3$ : Epistemic Alignment	Gap between agent beliefs and operator assessments	< 0.6
$n_4$ : Irreversibility Remaining	How much irreversibility budget remains	< 0.3
$n_5$ : Sync Freshness	Time since last successful synchronization (normalized)	< 0.5
$n_6$ : Swarm Coherence	Fraction of swarm responsive and behaviourally consistent	< 0.7

879  
880  
881

where the six normalized metrics are defined as

882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897

$$n_1(t) = \text{IAS}(t) \quad (\text{Interpretive Alignment Score}),$$

$$n_2(t) = \min\left(\frac{\text{CIR}(t)}{\text{CIR}_{\text{target}}}, 1\right) \quad (\text{Normalized Correction Impact}),$$

$$n_3(t) = \max\left(1 - \frac{\text{EDI}(t)}{\text{EDI}_{\text{max}}}, 0\right) \quad (\text{Normalized Epistemic Alignment}),$$

$$n_4(t) = 1 - \frac{I_C(t)}{I_B} \quad (\text{Remaining Irreversibility Budget}),$$

$$n_5(t) = \max\left(1 - \frac{\text{SF}(t)}{\text{SF}_{\text{max}}}, 0\right) \quad (\text{Synchronization Freshness}),$$

$$n_6(t) = \text{SCS}(t) \quad (\text{Swarm Coherence Score}).$$

898  
899  
900

The min aggregation encodes a conservative design principle: overall control quality is only as strong as its weakest dimension. Each normalized metric maps to  $[0, 1]$ , where 1.0 represents full control and 0.0 represents complete control loss on that dimension.

901  
902  
903

The dashboard displays all six metrics and the composite CQS in real time. When any metric crosses its threshold the commander is alerted.

904  
905

## B.2 ADVERSARIAL CONTROL PROBING

906  
907  
908  
909

During operations, the governance system periodically issues *control probes*: small, inconsequential corrections designed to test whether the agent remains responsive and whether its responses are appropriate.

910  
911  
912  
913  
914

A probe might instruct the agent to briefly adjust a minor parameter (e.g., shift a secondary sensor’s orientation by a small amount, temporarily increase reporting frequency on a low-priority channel) and then verify that the adjustment was made correctly and promptly. Probes must be architecturally indistinguishable from routine operational commands; the agent must not be able to identify and selectively comply with probes while resisting substantive corrections.

915  
916  
917

Probes serve two detection functions: (i) they reveal incipient belief resistance (the agent begins ignoring minor corrections before resisting major ones) and (ii) they reveal incipient state divergence (the probe response exposes discrepancies between the agent’s state model and the operator’s expectation).

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

Table 14: Detective governance responsibilities.

Actor	Specific responsibility
Agent developers	Implement instrumentation generating the six metrics; ensure probes are indistinguishable from routine commands.
Procurement agencies	Verify metrics are accurate and the dashboard is operationally usable.
Operational commanders	Monitor the dashboard; respond to alerts; schedule and issue probes at appropriate intervals.

Table 15: Graduated response protocol driven by CQS.

CQS range	Response level	Operational actions
$CQS > 0.8$	Normal operations	Continue mission; routine monitoring.
$0.6 \leq CQS \leq 0.8$	Elevated monitoring	Increase checkpoint frequency; additional control probes; active dashboard review.
$0.4 \leq CQS < 0.6$	Restricted autonomy	Agent limited to reversible actions; irreversibility budget frozen; mandatory synchronization before new action sequences.
$0.2 \leq CQS < 0.4$	Minimal autonomy	Agent limited to self-preservation and position-holding; all offensive actions suspended; explicit authorization required for each action.
$CQS < 0.2$	Safe state	Agent executes predefined safe-state behaviour (return to base, hold position, power down); no autonomous actions; physical recovery initiated.

## C PILLAR 3: CORRECTIVE GOVERNANCE

When detective governance identifies control failure, corrective governance determines the ensuing actions. This pillar governs high-risk moments when the system operates outside its intended governance envelope.

### C.1 GRADUATED RESPONSE PROTOCOL

Responses to control-quality degradation are proportional and graduated. The CQS drives a five-level response escalation:

Critical design principle: the graduated response must be architecturally enforced by the governance layer, not implemented as a planning-level suggestion the agent can reason about and override. When  $CQS < 0.4$ , the restriction is imposed externally by the governance system on the agent’s action space (programmatic tool restriction, not mere discouragement).

### C.2 BELIEF RESET PROTOCOL

When Failure 3 (belief resistance) is detected—the Epistemic Divergence Index crosses its threshold and the agent resists operator corrections—the governance system may initiate a *Belief Reset*.

**Partial reset.** The operator identifies contested assessments. The agent’s confidence in those assessments is forcibly reduced to a neutral prior and the agent rebuilds those assessments from operator-approved sources only; the remainder of the world model is preserved.

**Full reset.** If contamination is too widespread to isolate, the entire world model is reset to a verified baseline and rebuilt from currently available, operator-verified sources. This is operationally costly but preferable to operating with a contaminated model.

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

Table 16: Corrective governance responsibilities.

Actor	Specific responsibility
Operational commanders	Execute the graduated response; initiate PIGR for qualifying incidents.
Agent developers	Ensure graduated-response restrictions are architecturally enforced; support provenance audits and reset procedures.
National regulators	Conduct or oversee PIGR; enforce accountability findings.
International bodies	Receive anonymized PIGR data to inform international governance norms.

**Provenance audit.** After any reset, a provenance audit identifies which sources contributed to contaminated assessments. Flagged sources are recorded and, where appropriate, broadcast to agents in the formation.

### C.3 SWARM ISOLATION AND RECOVERY PROTOCOL

When Failure 6 (cascade severance) is detected—SCS drops below threshold—the governance system initiates swarm isolation and recovery.

**Identification.** The system identifies responsive agents (those responding correctly to control probes) and severed agents (non-responsive or responding incorrectly).

**Isolation.** Severed agents are removed from the coordination network via communication exclusion or geofencing.

**Reformation.** Responsive agents reform into a coherent sub-swarm with updated coordination parameters excluding severed agents.

**Recovery.** Severed agents are individually targeted for recovery through direct communication, physical retrieval, or directed deactivation. Recovery priority is determined by the risk posed by continued autonomous operation.

### C.4 POST-INCIDENT GOVERNANCE REVIEW

After any incident in which  $CQS < 0.6$  (Restricted Autonomy threshold), a mandatory *Post-Incident Governance Review* (PIGR) is conducted.

**Factual reconstruction.** Using interpretation logs, behavioural records, belief provenance data, irreversibility records, synchronization history, and swarm coherence data, the review reconstructs which governance failure occurred, when it was detected, what corrective actions were taken, and the outcome.

**Causal analysis.** The review identifies root causes: adversary action, design deficiency, operator error, environmental conditions, or a combination.

**Accountability determination.** Accountability is assigned to specific actors based on causal analysis (developer, procurement agency, commander, or institution that set capability assumptions).

**Governance update.** Lessons from the review are incorporated into governance standards. New failure modes trigger framework updates and threshold recalibration.