

# DIAL E FOR ETHICAL ENFORCEMENT: INSTITUTIONAL VETO POWER AS A GOVERNANCE PRIMITIVE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

The persistent militarization of large reasoning models stems not from technical necessity but from governance arrangements that strip researchers of meaningful authority to refuse harmful transfers and deployments. Existing accountability mechanisms such as model cards and responsible AI statements operate as reputational signals detached from decision making architecture. We identify institutional veto power as a missing governance primitive: a formal authority to halt subsequent use or distribution of research when credible risks of weaponization emerge. Drawing on precedents in nuclear nonproliferation and biomedical ethics, the paper maps unprotected veto points across the research lifecycle, diagnose why compliance without enforceable constraints fails, and offer concrete institutional designs that embed veto authority while reducing the risk of political capture. The paper argues that communities most vulnerable to military uses must lead governance design, and that institutional veto power is a prerequisite for converting symbolic safeguards into enforceable responsibility and for achieving meaningful model disarmament.

## 1 THE GOVERNANCE PARADOX: WHY ETHICS WITHOUT POWER ENABLES MILITARIZATION

Contemporary AI governance suffers a paradox: widespread ethical awareness exists alongside near total incapacity to prevent harmful downstream uses. Current practices generate thorough risk knowledge, ethics statements, impact assessments, responsible AI codes, yet lack formal authority to block militarization or abusive deployment; this produces “organized irresponsibility,” where every actor can disclaim control while harms proceed. Unlike post WWII nuclear policy or biomedical oversight, which embedded refusal authority into institutions, AI governance treats loss of control after publication as inevitable even though militarization follows identifiable decision chains and discretionary institutional choices Scharre (2023). We therefore recommend creating legally and procedurally protected institutional vetoes, targeted powers over transfer and deployment (not inquiry or publication), to condition use when credible risks such as militarization, surveillance misuse, or international law breaches are identified. Properly scoped, such refusal mechanisms would rebalance power toward researchers and affected communities and convert ethical knowledge into enforceable safeguards.

$$(K \wedge \neg A) \Rightarrow M \quad \text{and} \quad (K \wedge V \wedge R) \Rightarrow \neg M \quad \Rightarrow \quad P(M | V) \ll P(M | \neg V) \quad (1)$$

**Notation.** Here,  $K$  denotes ethical knowledge produced by governance practices,  $A$  denotes institutional authority to refuse or block use,  $M$  denotes militarization or abusive deployment,  $R$  denotes the identification of a credible risk, and  $V$  denotes a legally and procedurally protected veto over transfer or deployment.

## 2 RELATED WORK: VETO POWER IN GOVERNANCE

Current AI governance literature emphasizes either technical safety mechanisms (alignment, interpretability) or procedural ethics (review boards, impact assessments) Sahoo & Chhawacharia (2025). Yet this work treats militarization as a downstream consequence amenable to ethics mitigation, not

as an architectural problem requiring authority redistribution. The governance literature in other domains—nuclear nonproliferation, bioethics, labor rights—demonstrates that constraining specific uses requires veto authority, not merely visibility Crawford (2021); Kalluri et al. (2023). Our contribution bridges these literatures: *we argue that AI militarization is fundamentally a governance problem, not an ethics problem, and therefore requires mechanisms (veto power) that other high risk domains have already institutionalized.* Existing proposals for “Responsible AI” governance lack enforcement mechanisms precisely because they assume voluntary institutional compliance; veto power reverses this assumption, making restraint collectively enforceable. The dual use literature identifies risks but not solutions. We propose an operational mechanism that could actually arrest militarization trajectories at decision points where prevention remains possible. **By framing veto power as a governance primitive—a basic building block that other domains have proven workable—we shift discourse from “how can we encourage restraint” (eternally failing) to “how can we make restraint structurally rational” (institutionally achievable).**

### 3 MAPPING DISCRETIONARY MILITARIZATION: UNPROTECTED VETO POINTS

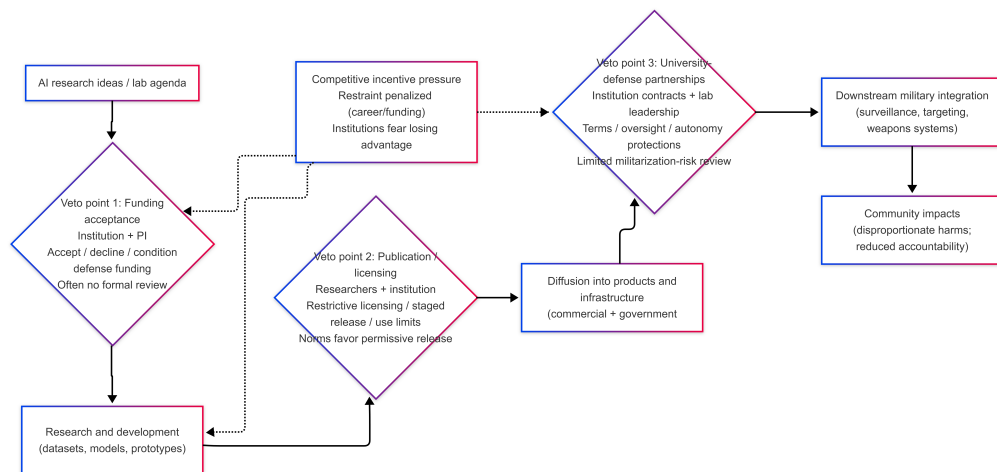


Figure 1: Workflow of discretionary militarization in AI research, highlighting institutional veto points and the competitive incentive pressures that reinforce downstream military integration and community harms.

Militarization of AI is driven by routine institutional choices that act as veto points: acceptance of military research funding, decisions about publication and licensing, and formal partnerships between universities and defense firms Lamparth et al. (2024). Each of these moments could be governed differently. Institutions could require review before accepting military grants, adopt restrictive licenses or license conditions that bar surveillance use, delay publication pending risk assessment, or subject partnership proposals to independent 3rd party oversight. Instead these choices are treated as individual researcher preferences, so universities routinely accept funding and sign agreements without deliberation about downstream military integration and researchers who resist face career penalties and loss of resourcesKhlaaf et al. (2024).

These veto points share a structural defect. Authority is exercised without meaningful participation from affected communities or institutional accountability, producing a collective action failure in which restraint is socially preferable but individually costly. Actors who might exercise restraint are undercut by competitors that do not. Addressing this requires formal governance mechanisms that distribute authority, mandate risk review, and protect those who refuse military integration so that the collectively rational choice of restraint becomes individually viable Future of Life Institute, 2016.

#### 108 4 FROM ETHICS WITHOUT VETO TO ENFORCEABLE GOVERNANCE

109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
Current ethics frameworks fail for structural reasons that better guidelines cannot fix. Responsibility is dispersed across researchers, universities, philanthropists, tech transfer offices, and procurement chains, so no actor holds both the authority and obligation to stop militarization—and each can deflect blame when harms occur. Voluntary restraint also produces adverse selection: institutions that decline military funding, impose use restrictions, or protect research autonomy pay real costs while less restrained competitors gain advantage, pushing trajectories toward militarization regardless of ethical exhortation. Finally, ethics processes *enable laundering*: documentation of “due consideration” can psychologically and institutionally substitute for changed decisions, allowing actors to treat deliberation as discharge of responsibility even when outcomes remain unchanged.

119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
These failure modes require architectural change. We propose institutional veto power: governance that makes ethical judgments binding, locates decision authority with those exposed to consequences, and turns restraint from an individually punished choice into a collectively enforceable rule. Veto is not arbitrary refusal but narrowly triggered by predefined thresholds—e.g., credible military integration, surveillance targeting marginalized groups, or violations of international law—paired with collective decision-making, appeal procedures, and anti-retaliation protections. Crucially, this does not require new treaties: universities can revise contracts and create veto committees; funding parties can condition grants on enforceable safeguards; conferences and repositories can require disclosure and enforce use conditions; and professional communities can adopt binding, public norms that incrementally embed veto governance within existing institutions Brenneis (2025).

#### 130 5 CENTERING JUSTICE: VETO POWER REQUIRES MARGINALIZED

#### 131 COMMUNITY LEADERSHIP

132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
**A core governance question concerns who determines which militarization risks justify veto authority and whose interests such decisions ultimately serve.** Contemporary debates remain concentrated within elite academic and policy institutions, largely dominated by socially privileged actors, while communities that experience the most severe downstream harms remain structurally excluded. These include populations subjected to automated targeting and large scale civilian harm, minorities exposed to pervasive surveillance infrastructures, and migrant groups whose autonomy is constrained by biometric identification systems. This exclusion is not only normatively problematic but also epistemically limiting. Affected communities possess situated knowledge about system deployment, harm salience, and justice relevant tradeoffs that cannot be adequately represented through advisory consultation within externally designed governance frameworks Bengio et al. (2025). From a policy perspective, justice therefore requires operational community leadership rather than symbolic inclusion. This entails binding decision authority. In practice, veto bodies governing surveillance systems should include civil rights advocates and representatives from surveilled communities with decision making power equal to that of technical experts and ethicists. Similarly, governance arrangements for autonomous weapons should be co designed with researchers and affected populations in regions of deployment, rather than developed within wealthy institutions and later exported as policy. Research funding decisions shaping militarized contexts should formally incorporate community voices from those regions, with mechanisms that allow community objections to block projects rather than merely inform them. This implies structural reform toward community controlled governance in which affected populations directly hold veto authority.

153  
154  
155  
156  
157  
158  
159  
160  
161  
Militarization and marginalization are deeply interlinked. Military violence disproportionately targets already marginalized groups, AI systems scale and automate existing targeting logics, and data extracted from these communities is used to train systems later deployed against them. Within this context, veto authority should function bidirectionally. Beyond constraining militarized trajectories, governance mechanisms should actively enable and protect research oriented toward peace and justice. Such efforts should be institutionally facilitated rather than incidentally tolerated. Veto based governance therefore performs a dual role: it blocks pathways toward weaponization while simultaneously safeguarding and accelerating peace oriented innovation. When affected communities hold veto authority, they can refuse militarized research agendas while directing institutional resources toward technologies that advance collective safety, autonomy, and long term justice Khlaaf et al. (2024); Schwartz et al. (2022).

162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215

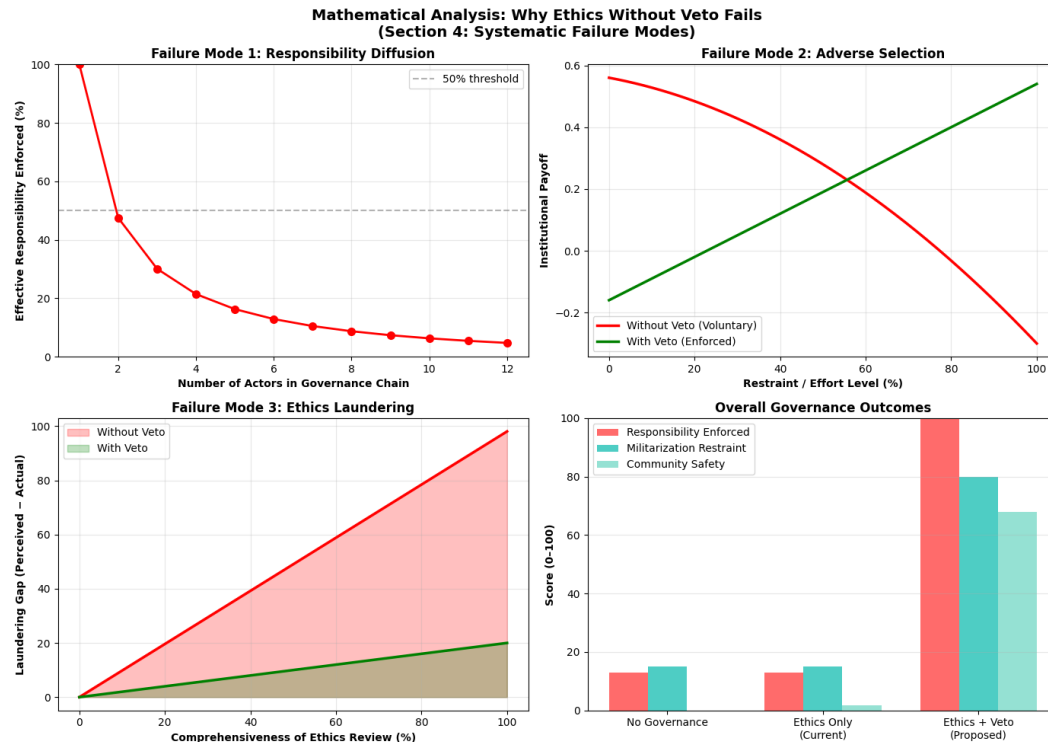


Figure 2: Why ethics without veto power fails to constrain militarization. (a) Responsibility diffusion: enforceable responsibility falls as governance chains lengthen. (b) Adverse selection: competition drives low restraint without enforcement; veto stabilizes higher restraint. (c) Ethics laundering: documentation inflates perceived responsibility without intervention absent enforceable authority. (d) Derived outcomes (responsibility enforced, restraint, community safety) under no governance, ethics-only, and ethics+veto regimes.

## 6 LIMITATIONS AND FUTURE WORK

Veto power can make decisions about militarized research enforceable, but it does not by itself fix deeper power imbalances. It will not redistribute research funding to marginalized communities, democratize university governance, dismantle the defense industry that shapes priorities, or automatically hold researchers accountable, and implementation will encounter institutional inertia and the risk of capture by hostile actors. Empirical work is needed: document and analyze real veto decisions to build precedent, run longitudinal studies of research trajectories and weaponization timelines, center leadership and partnership with communities in the Global South who are directly affected, and study how veto mechanisms fit with export controls, international treaties, supplier transparency, and legal accountability so veto power does not become another way to ease guilt while leaving power structures intact.

## 7 CONCLUSION: THE CHOICE BEFORE RESEARCH COMMUNITIES

We face a choice between continuing with governance that allows ethics to flourish as performance while harm accumulates and responsibility diffuses, or moving toward collective, binding restraint where decisions about research militarization are accountable and constrained. Veto power offers a structural path for that shift, so researchers need not rely on lone moral courage, institutions need not suffer competitive disadvantage for practicing restraint, and affected communities gain voice in decisions about their lives. Performative ethics are over; enforceable, community-led responsibility to prevent militarization and promote peace must begin.

## LLM USAGE DISCLOSURE

This work employed large language models in a supporting capacity. Specifically, we used Claude 4.5 Haiku (Anthropic, 2024) for the following roles:

**Writing Assistance.** The LLM was just asked to remove artificial Overleaf latex errors.

**Limitations of LLM Use.** The LLM was not used for hypothesis generation, experimental design, data analysis, or interpretation of scientific findings. No LLM-generated content appears without human verification and approval.

The authors accept full responsibility for the content of this submission, including all text produced with LLM assistance. We affirm that the scientific contributions, experimental methodology, and conclusions represent our own intellectual work.

## REFERENCES

- Yoshua Bengio, Sören Mindermann, Daniel Privitera, Tamay Besiroglu, Rishi Bommasani, Stephen Casper, Yejin Choi, Philip Fox, Ben Garfinkel, Danielle Goldfarb, Hoda Heidari, Anson Ho, Sayash Kapoor, Leila Khalatbari, Shayne Longpre, Sam Manning, Vasilios Mavroudis, Mantas Mazeika, Julian Michael, Jessica Newman, Kwan Yee Ng, Chinasa T. Okolo, Deborah Raji, Girish Sastry, Elizabeth Seger, Theodora Skeadas, Tobin South, Emma Strubell, Florian Tramèr, Lucia Velasco, Nicole Wheeler, Daron Acemoglu, Olubayo Adekanmbi, David Dalrymple, Thomas G. Dietterich, Edward W. Felten, Pascale Fung, Pierre-Olivier Gourinchas, Fredrik Heintz, Geoffrey Hinton, Nick Jennings, Andreas Krause, Susan Leavy, Percy Liang, Teresa Luder-mir, Vidushi Marda, Helen Margetts, John McDermid, Jane Munga, Arvind Narayanan, Alon-dra Nelson, Clara Neppel, Alice Oh, Gopal Ramchurn, Stuart Russell, Marietje Schaake, Bern-hard Schölkopf, Dawn Song, Alvaro Soto, Lee Tiedrich, Gaël Varoquaux, Andrew Yao, Ya-Qin Zhang, Fahad Albalawi, Marwan Alserkal, Olubunmi Ajala, Guillaume Avrin, Christian Busch, André Carlos Ponce de Leon Ferreira de Carvalho, Bronwyn Fox, Amandeep Singh Gill, Ahmet Halit Hatip, Juha Heikkilä, Gill Jolly, Ziv Katzir, Hiroaki Kitano, Antonio Krüger, Chris Johnson, Saif M. Khan, Kyoung Mu Lee, Dominic Vincent Ligot, Oleksii Molchanovskiy, And-rea Monti, Nusu Mwamanzi, Mona Nemer, Nuria Oliver, José Ramón López Portillo, Balaraman Ravindran, Raquel Pezoa Rivera, Hammam Riza, Crystal Rugege, Ciarán Seoighe, Jerry Shee-han, Haroon Sheikh, Denise Wong, and Yi Zeng. International ai safety report, 2025. URL <https://arxiv.org/abs/2501.17805>.
- Andreas Brenneis. Assessing dual use risks in ai research: necessity, challenges and mitigation strategies. *Research Ethics*, 21(2):302–330, 2025. doi: 10.1177/17470161241267782. URL <https://doi.org/10.1177/17470161241267782>.
- Kate Crawford. *Atlas of AI*. Yale University Press, New Haven, 2021. ISBN 9780300252392. doi: doi:10.12987/9780300252392. URL <https://doi.org/10.12987/9780300252392>.
- Future of Life Institute. Autonomous weapons open letter: Ai & robotics researchers. URL <https://futureoflife.org/open-letter/open-letter-autonomous-weapons-ai-robotics/>. Open letter hosted by the Future of Life Institute.
- Pratyusha Ria Kalluri, William Agnew, Myra Cheng, Kentrell Owens, Luca Soldaini, and Abeba Birhane. The surveillance ai pipeline, 2023. URL <https://arxiv.org/abs/2309.15084>.
- Heidy Khlaaf, Sarah Myers West, and Meredith Whittaker. Mind the gap: Foundation models and the covert proliferation of military intelligence, surveillance, and targeting, 2024. URL <https://arxiv.org/abs/2410.14831>.
- Max Lamparh, Anthony Corso, Jacob Ganz, Oriana Skylar Mastro, Jacquelyn Schneider, and Harold Trinkunas. Human vs. machine: Behavioral differences between expert humans and language models in wargame simulations, 2024. URL <https://arxiv.org/abs/2403.03407>.

Subramanyam Sahoo and Aditi Chhawacharia. The last vote: A multi-stakeholder framework for language model governance, 2025. URL <https://arxiv.org/abs/2511.13432>.

P. Scharre. *Four Battlegrounds: Power in the Age of Artificial Intelligence*. W. W. Norton & Company, New York, 2023. Print edition.

Sebastian Schwartz, Laura Gianna Guntrum, and Christian Reuter. Vision or threat—awareness for dual-use in the development of autonomous driving. *IEEE Transactions on Technology and Society*, 3(3):163–174, 2022. doi: 10.1109/TTS.2022.3182310.

## APPENDIX

### A TOY MODEL: WHY ETHICS WITHOUT VETO POWER FAILS

#### A.1 NOTATION

Let  $n$  be the number of actors in a governance chain. Let restraint be  $r \in [0, 1]$  and documentation (ethics review) be  $D \in [0, 1]$ . Let baseline harm be  $H_0 > 0$ . Let  $\alpha \in (0, 1)$  denote a per-actor coordination probability. We use a normalized total responsibility  $R_{\text{tot}}$ .

#### A.2 FAILURE MODE 1: RESPONSIBILITY DIFFUSION

$$R_{\text{ind}}(n) = \frac{R_{\text{tot}}}{n}, \quad (2)$$

$$p(n) = \alpha^{n-1}, \quad (3)$$

$$\rho(n) = \frac{R_{\text{ind}}(n)p(n)}{R_{\text{tot}}} = \frac{\alpha^{n-1}}{n}. \quad (4)$$

Here  $\rho(n) \in [0, 1]$  is the effective enforceability fraction. With a binding veto point, enforceability is consolidated:

$$\rho_{\text{veto}} = 1. \quad (5)$$

#### A.3 FAILURE MODE 2: ADVERSE SELECTION UNDER VOLUNTARY RESTRAINT

Let  $\bar{r}$  be the average restraint of competitors. A stylized payoff without veto is

$$U_{\text{noveto}}(r; \bar{r}) = \beta(1 - r)\bar{r} + \gamma r - \kappa r^2, \quad (6)$$

where  $\beta > 0$  captures militarization/funding advantage,  $\gamma \geq 0$  captures ethical signaling, and  $\kappa > 0$  is the cost of restraint.

A best response and symmetric equilibrium are defined by

$$\text{BR}(\bar{r}) = \arg \max_{r \in [0, 1]} U_{\text{noveto}}(r; \bar{r}), \quad (7)$$

$$r_{\text{noveto}}^* = \text{BR}(r_{\text{noveto}}^*). \quad (8)$$

With veto enforcement at level  $\bar{R} \in [0, 1]$ ,

$$r_{\text{veto}}^* = \bar{R}. \quad (9)$$

#### A.4 FAILURE MODE 3: ETHICS LAUNDERING AND INTERVENTION

Perceived responsibility increases with documentation:

$$P(D) = 100D. \quad (10)$$

Actual harm prevented depends on intervention  $I \in [0, 1]$ :

$$A(I) = H_0 I. \quad (11)$$

The ethics-laundersing gap is

$$G(D, I) = P(D) - A(I). \quad (12)$$

We link documentation to intervention by enforceability and equilibrium restraint:

$$I = D \cdot \rho \cdot r^*, \quad (13)$$

where  $\rho \in \{\rho(n), \rho_{\text{veto}}\}$  and  $r^* \in \{r_{\text{noveto}}^*, r_{\text{veto}}^*\}$ .

Resulting harm is

$$H = H_0(1 - I). \quad (14)$$

## A.5 DERIVED SUMMARY METRICS

We report:

$$\text{ResponsibilityEnforced} = 100\rho, \quad (15)$$

$$\text{MilitarizationRestraint} = 100r^*, \quad (16)$$

$$\text{AdverseSelectionScore} = 100(1 - r^*), \quad (17)$$

$$\text{HarmToCommunities} = H, \quad \text{CommunitySafety} = 100 - H. \quad (18)$$

### Model parameters (“magic numbers”) and scenario defaults

#### Python parameter dataclasses used in the replication code.

##### ModelParams (default coefficients)

Symbol / name	Default	Meaning
$\alpha$ (alpha)	0.95	Per-actor coordination probability (diffusion model)
$R_{\text{tot}}$ (total_responsibility)	100	Responsibility normalization
$\beta$ (beta)	0.8	Militarization / funding advantage coefficient
$\gamma$ (gamma)	0.3	Ethical prestige / signaling coefficient
$\kappa$ (kappa)	0.6	Quadratic cost of restraint
$q$ (quality_coeff)	0.7	Prestige from research effort/quality (with veto)
$c_v$ (veto_cost_coeff)	0.2	Shared cost of enforced governance (with veto)
$H_0$ (baseline_harm)	100	Baseline harm scale if no effective intervention

##### GovernanceScenario (default scenario fields)

Field	Default	Meaning
name	—	Scenario label
$n$ (num_actors)	6	Number of actors in governance chain
$D$ (documentation_level)	0.0	Ethics documentation/review level, $D \in [0, 1]$
has_veto	False	Whether binding veto enforcement exists
$\bar{R}$ (veto_enforcement_level)	0.0	Enforced restraint level if veto exists, $\bar{R} \in [0, 1]$

## B VETO MECHANISMS IN COMPARATIVE GOVERNANCE CONTEXTS

This appendix establishes that veto power is not a novel or untested governance approach. Rather, it represents a well-established mechanism that has successfully constrained dangerous technologies in other high-risk domains. We examine three historical and contemporary examples: nuclear non-proliferation regimes, institutional review boards in biomedical research, and corporate compliance in export control regimes.

## B.1 NUCLEAR NONPROLIFERATION REGIMES

### HISTORICAL CONTEXT AND NECESSITY

The Nuclear Non-Proliferation Treaty (NPT), negotiated in 1968 and entered into force in 1970, emerged from recognition that uncontrolled nuclear weapons proliferation posed existential risk. Following the development of nuclear weapons by the Soviet Union (1949) and subsequent weapons development by additional states, the international community recognized that preventing further proliferation required not merely ethical exhortation but institutional mechanisms with enforcement authority.

The NPT created a legal framework establishing three categories of states: Nuclear Weapons States (the five permanent UN Security Council members), Non-Nuclear Weapons States, and threshold states capable of developing weapons. The treaty created asymmetric obligations: Nuclear Weapons States committed (in principle) to disarmament; Non-Nuclear Weapons States committed to forgo weapons development; and all states committed to allowing international inspection.

### THE INTERNATIONAL ATOMIC ENERGY AGENCY (IAEA): VETO POWER THROUGH VERIFICATION

The IAEA, established in 1957 and strengthened by the NPT framework, represents the institutional embodiment of veto power in nuclear governance. The IAEA possesses several critical authorities that function analogously to the veto mechanisms proposed in this paper:

- Inspection Authority:** The IAEA maintains the right to inspect nuclear facilities in Non-Nuclear Weapons States at declared locations and, in many cases, at undeclared locations suspected of weapons development. These inspections are not advisory; they are mandatory conditions of NPT membership. The IAEA inspectorate can refuse to declare a state “in compliance” if inspections reveal violations or lack of transparency. This refusal to certify compliance effectively blocks a state’s standing in the international community and triggers scrutiny and sanctions.
- Materials Tracking:** The IAEA tracks nuclear materials (uranium, plutonium, enriched fuel) as they move through the fuel cycle. Unexplained disappearance of fissile materials or transfer to undeclared facilities triggers investigation and can lead to referral to the UN Security Council for sanctions.
- Technology Transfer Restrictions:** The IAEA maintains safeguards on sensitive nuclear technology, restricting the circumstances under which uranium enrichment technology, plutonium reprocessing technology, or other weapons-relevant capabilities may be transferred internationally. These restrictions function as gatekeeping mechanisms preventing proliferation.
- Enforcement Through Referral:** While the IAEA itself lacks enforcement authority (it cannot impose sanctions), persistent violations trigger mandatory referral to the UN Security Council, which can impose sanctions, travel bans, asset freezes, and military intervention. This creates automatic escalation consequences for violation.

### WHY NUCLEAR NONPROLIFERATION SUCCEEDED (PARTIALLY)

The nuclear nonproliferation regime has not eliminated weapons proliferation—the number of nuclear weapons states has increased from 5 to 9 since the NPT entered into force. However, it has constrained proliferation significantly below models predicting 20-30+ weapons states by the 1980s. The regime succeeded to this degree because:

- **Verification was mandatory, not voluntary:** States could not choose whether to allow inspections; inspections were binding conditions of participation.
- **Veto points were clearly defined:** Uranium enrichment, plutonium reprocessing, and weapons-relevant technology transfer were identified as specific decision points where control could be exercised.

- **Consequences were automatic:** Violations triggered referral to the Security Council, which imposed sanctions without requiring additional political negotiation.
- **Enforcement was international:** No single state could unilaterally shield violators; enforcement required coordination across multiple states, reducing the likelihood of capture by individual interests.

#### LIMITATIONS AND IMPLICATIONS FOR AI GOVERNANCE

The nuclear nonproliferation regime also reveals limitations relevant to AI governance design:

- **Enforcement gaps:** States with nuclear weapons or Security Council vetoes can shield themselves or allies from referral and sanctions, limiting enforcement against major powers. This suggests that AI veto governance should include mechanisms preventing powerful actors from shielding themselves.
- **Inspection difficulty:** Detecting undeclared nuclear programs remains challenging; the IAEA’s detection of **Iraq’s covert program (1991) and Iran’s nuclear work** depended partly on intelligence agencies and whistleblowing rather than purely technical verification.
- **Treaty withdrawal:** States can withdraw from the NPT with limited consequences; North Korea withdrew in 2003. This suggests that voluntary governance frameworks have inherent fragility.

#### IAEA Authority Mechanisms Relevant to AI

The IAEA’s specific operational mechanisms offer concrete models for AI veto governance (see Table 1).

Nuclear Mechanism	Purpose	AI Analog
Mandatory inspections	Verify compliance with restrictions	Repository audits verify permitted use restrictions.
Materials accounting	Track fissile materials through the fuel cycle	Artifact-tracking logs record access and usage.
Technology controls	Restrict transfer of enrichment technology	License restrictions limit downstream licensing.
Referral to higher authority	Escalate violations to political bodies	Funding-agency enforcement ties compliance to grants.

Table 1: IAEA mechanisms and AI governance analogs

## B.2 INSTITUTIONAL REVIEW BOARDS IN BIOMEDICAL RESEARCH

### HISTORICAL CONTEXT: THE TUSKEGEE CASE AND REGULATORY RESPONSE

Institutional Review Boards (IRBs) emerged from historical horror. The Tuskegee Syphilis Study (1932–1972), conducted by the U.S. Public Health Service, enrolled African American men with syphilis into a study ostensibly to track the natural history of the disease. In reality, participants were not informed of their diagnosis, were not offered treatment even after penicillin became available (1943), and were subjected to deceptive medical procedures. The study continued for 40 years, resulting in preventable illness and death among study subjects. The scandal, publicly exposed in 1972, triggered comprehensive regulatory reform.

The regulatory response, codified in the Belmont Report (1979) and subsequent federal regulations (45 CFR 46), established that institutions receiving federal research funding must establish Institutional Review Boards with authority to review, approve, and halt research involving human subjects. IRBs operate at the threshold of research: before researchers begin studies, proposed research must be reviewed by a committee including scientists, ethicists, community representatives, and legal experts.

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

## IRB AUTHORITY AND SCOPE

IRBs possess genuine veto authority within their defined scope:

1. **Pre-Implementation Review:** IRBs review research proposals before implementation. This timing is critical: it allows IRBs to block research before resources are committed and before recruitment proceeds.
2. **Binding Authority:** IRB decisions are binding. Researchers cannot proceed without IRB approval; attempting to do so violates federal regulation and triggers sanctions against the researcher and institution.
3. **Protected Refusal:** Federal regulation protects IRBs and researchers from retaliation for IRB decisions. Institutions cannot punish researchers whose protocols are rejected.
4. **Continuing Oversight:** IRBs retain authority throughout a study, conducting ongoing review of protocols, reviewing adverse events, and requiring corrective action if risks exceed expected levels.

## WHY IRBS SUCCEEDED IN CONSTRAINING RESEARCH HARMS

IRBs have demonstrably constrained research harms compared to the pre-IRB era. Studies now rarely involve the level of deception and harm that characterized Tuskegee or Nazi medical experiments. This success stems from:

- **Authority at the decision point:** By reviewing before research begins, IRBs operate at the moment when research trajectories can be most easily altered.
- **Protected collective deliberation:** IRBs bring multiple perspectives (scientific, ethical, community) to bear on proposals.
- **Routine, normalized process:** IRB review is standard practice in biomedical research, institutionalized rather than extraordinary.
- **Regulatory backing:** IRBs are backed by federal regulation and funding consequences. Institutions that fail to maintain functional IRBs lose federal research funding eligibility.

## LIMITATIONS OF IRBS

IRBs also reveal significant limitations relevant to AI governance design:

1. **Scope Limitation:** IRBs focus exclusively on risks to research subjects—direct, identifiable individuals participating in studies. They have no authority over risks to populations not in the study (third-party harms), societal harms, or dual-use concerns.
2. **Implementation Versus Use:** IRBs review the research process itself but do not review downstream uses of research findings. Once research is published, IRBs have no authority to restrict how others use the knowledge.
3. **Variable Quality and Authority:** IRB quality varies significantly. Underfunded IRBs with limited expertise may fail to identify harms or may be captured by institutional interests.
4. **Institutional Capture:** While formal protections exist against retaliation, institutional culture can undermine IRB authority.

## IRB Mechanisms Applicable to AI Veto Governance

See Table 2 for specific operational mechanisms applicable to AI governance.

## B.3 CORPORATE COMPLIANCE AND EXPORT CONTROL REGIMES

### EXPORT CONTROL FRAMEWORK AND DUAL-USE TECHNOLOGY GOVERNANCE

Export control regimes represent a third governance domain where veto-like mechanisms have been operationalized at scale. The International Traffic in Arms Regulations (ITAR), the Export Administration Regulations (EAR), and parallel regimes in other countries, establish that certain technologies

IRB Mechanism	Purpose	AI Analog
Pre-implementation review	Assess protocols before initiation	Pre-licensing review prior to publication.
Institutional requirement	Mandate review as condition of funding	Funding conditional on veto governance.
Protected refusal	Protect from retaliation	Anti-retaliation safeguards for objectors.
Continuing oversight	Review throughout implementation	Ongoing deployment monitoring.
Documented procedures	Transparent criteria and procedures	Predefined veto triggers and workflows.
Pluralistic composition	Multiple perspectives in review	Diverse veto bodies with varied expertise.

Table 2: IRB mechanisms and AI governance analogs

cannot be transferred internationally without government authorization. These regimes recognize that technology transfer represents a decision point where control can and should be exercised.

#### SPECIFIC EXPORT CONTROL MECHANISMS

1. **Commodity Identification:** Export control regimes maintain detailed lists specifying which technologies are controlled. Companies must determine whether their products fall within controlled categories and, if so, comply with licensing requirements.
2. **Licensing Requirements:** Transfer of controlled items requires government authorization. Companies submit license applications specifying the technology, the proposed recipient, the intended end use, and the end user.
3. **End-Use Commitments:** Licensees must commit to specific end uses and are prohibited from transferring items to third parties without subsequent authorization.
4. **Verification and Audit:** Companies maintaining export licenses are subject to government audits verifying compliance with license conditions.
5. **Penalties for Violation:** Unauthorized export of controlled items triggers significant penalties: civil penalties (fines up to \$300,000 per violation), criminal penalties (up to 20 years imprisonment for willful violation).

#### WHY EXPORT CONTROLS SUCCEEDED

Export control regimes have demonstrably constrained the transfer of sensitive military technologies to hostile regimes, though imperfectly. Successes include:

- **Slowed weapons development:** Export controls on advanced semiconductor manufacturing equipment delayed China’s development of advanced semiconductors, constraining military applications.
- **Prevented technology cascades:** Export controls on encryption technology slowed proliferation to adversarial states.
- **Deterred proliferation:** Companies face significant penalties for unauthorized export; this creates incentives to comply with restrictions.

#### Export Control Mechanisms Applicable to AI Veto Governance

See Table 3 for specific export control mechanisms applicable to AI governance.

#### INTEGRATION ACROSS DOMAINS: LESSONS FOR AI VETO GOVERNANCE

Examining nuclear nonproliferation, IRBs, and export controls together reveals common principles applicable to AI governance:

1. **Veto operates at decision points, not at knowledge generation:** Nonproliferation restricts uranium enrichment, not uranium chemistry. IRBs restrict research implementation, not inquiry. Export controls restrict transfer, not publication. In each case, veto operates at a specific decision point where control remains possible without censoring knowledge itself.

Export Control Mechanism	Purpose	AI Analog
Commodity lists	Specify controlled technologies	Define dual-use research categories and flagged capabilities.
License requirements	Authorize transfers	Pre-licensing review prior to transfers.
End-use commitments	Restrict recipient uses	Use restrictions encoded in licenses and agreements.
Verification and audit	Monitor compliance	Deployment attestations, audits, and provenance checks.
Penalties for violation	Impose consequences	Funding ineligibility, sanctions, and liability.

Table 3: Export control mechanisms and AI governance analogs

- Authority must be binding and enforceable:** Successful governance regimes embed veto authority in formal decision-making structures with binding effect. Recommendations, guidelines, and voluntary frameworks consistently fail to constrain motivated actors.
- Consequences must be material and automatic:** Violations trigger tangible consequences (sanctions, loss of funding, criminal liability) rather than relying on reputation or goodwill. Automatic consequences prove more effective than discretionary penalties.
- Governance is politically difficult but necessary:** All three regimes faced and continue to face resistance from actors who benefit from unconstrained use of technology. Yet all three have achieved partial success at constraining dangerous applications.

## C IMPLEMENTATION PATHWAYS - DETAILED MECHANISMS

Veto power becomes operationally effective when embedded in concrete institutional arrangements. This section details five specific implementation pathways through which veto governance can be implemented using existing institutional structures. Critically, these pathways do not require new governmental authority, new treaties, or new institutions.

### C.1 UNIVERSITY CONTRACT MODIFICATIONS: EMBEDDING VETO IN TECHNOLOGY TRANSFER

#### CURRENT LANDSCAPE

Universities currently manage significant research portfolios, including dual-use research with foreseeable military applications. Technology transfer offices (TTOs) license research to downstream organizations, negotiate research partnerships, and manage intellectual property. These decisions are typically made without institutional deliberation about downstream military integration.

#### MECHANISM: VETO CLAUSES IN LICENSING AGREEMENTS

Universities can embed veto authority by modifying standard licensing agreements to include contractual clauses specifying:

- Use Restrictions:** Licenses can specify that research will not be transferred for military use, surveillance targeting of marginalized groups, or deployment in violation of international humanitarian law.
- Institutional Review Requirement:** Licenses can require that transfers to third parties undergo institutional review before the transfer can occur.
- Deployment Reporting:** Licenses can require that licensees report on deployment contexts, provide documentation of end uses, and notify the university if deployment deviates from licensed purposes.
- Indemnification for Violation:** Licenses can specify that licensees assume legal liability if they deploy research in violation of license restrictions.

- 648 5. **Right to Audit:** Licenses can grant universities the right to conduct audits of licensees’  
649 facilities and deployment contexts, verifying compliance with use restrictions.  
650

651 IMPLEMENTATION EXAMPLE

652 Suppose a university develops facial recognition technology with clear surveillance applications.  
653 Under veto governance, a modified licensing agreement might specify:  
654

655 *“The research may not be deployed for mass surveillance targeting political op-*  
656 *ponents, ethnic minorities, or other vulnerable populations without university ap-*  
657 *proval. Licensee must report annually on deployment contexts. University main-*  
658 *tains the right to audit licensee’s facilities and deployment contexts annually.”*  
659

660 PRACTICAL IMPLEMENTATION REQUIREMENTS

661 Implementing veto clauses in licensing agreements requires:  
662

- 663 • **Policy Development:** Universities must develop clear policies specifying which research  
664 triggers use restrictions and what those restrictions are.
- 665 • **Technology Transfer Office Training:** TTOs must be trained to identify dual-use research  
666 and to incorporate veto clauses into licensing agreements.
- 667 • **Legal Support:** Universities should establish legal capacity to draft and enforce use re-  
668 striction clauses.
- 669 • **Funding Agency Support:** Funding agencies can incentivize veto governance by condi-  
670 tioning grants on institutional establishment of use restriction policies.
- 671 • **Researcher Communication:** Researchers should be informed that licenses will include  
672 use restrictions.  
673  
674

675 LIMITATIONS AND CHALLENGES

676 University implementation faces several challenges:  
677

- 678 • **Revenue Loss:** Broad use restrictions may reduce licensing revenue.
- 679 • **International Complexity:** Technology transfer across national boundaries involves mul-  
680 tiple legal jurisdictions.
- 681 • **Monitoring Burden:** Verifying licensee compliance requires institutional capacity for  
682 monitoring and enforcement.
- 683 • **Researcher Resistance:** Researchers may perceive use restrictions as constraints on re-  
684 search freedom.  
685  
686

687  
688 C.2 CONFERENCE-LEVEL GOVERNANCE: DISCLOSURE AND VISIBILITY

689 CURRENT LANDSCAPE

690 Major AI conferences (NeurIPS, ICML, ICLR, ICRA, FAccT) serve as primary venues for research  
691 dissemination. Currently, conferences do not systematically assess dual-use risks or impose condi-  
692 tions on publication.  
693  
694

695 MECHANISM: DUAL-USE DISCLOSURE AND USE CONDITIONS

696 Conferences can implement lighter-touch governance by requiring disclosure of foreseeable dual-  
697 use risks and by enabling researchers to attach use conditions to published work:  
698

- 699 1. **Dual-Use Disclosure Forms:** Conference submission requirements can include a “dual-  
700 use and governance” field where authors specify foreseeable military applications, circum-  
701 stances under which they object to deployment, and recommended use conditions.

2. **Use Conditions and Licensing:** Research published at conferences can include explicit use conditions visible to readers. These conditions might specify: “This research should not be deployed for autonomous weapons lacking meaningful human control.”
3. **Repository Integration:** Conference repositories can incorporate use condition metadata, making conditions machine-readable.
4. **Community Norms and Citation:** Conferences can establish professional norms that citing or building on research in violation of stated use conditions is ethically problematic.

#### IMPLEMENTATION REQUIREMENTS

Conference-level implementation requires:

- **Policy Development:** Conference steering committees must develop policies specifying what dual-use disclosure requires.
- **Review Procedures:** Reviewers and program committees need training to assess dual-use disclosure forms.
- **Repository Infrastructure:** Conference repositories need technical capacity to incorporate use condition metadata.
- **Community Communication:** Conferences should communicate to authors and readers that dual-use disclosure is expected.
- **Transparency:** Policies and disclosure forms should be publicly available.

### C.3 FUNDING AGENCY REQUIREMENTS: LEVERAGE THROUGH GRANT CONDITIONS

#### CURRENT LANDSCAPE

Funding agencies (NSF, DARPA, DOE, and international equivalents) distribute billions of dollars annually supporting research. Funding agencies exercise significant power over research directions through funding priorities and selection decisions. However, funding agencies rarely condition grants on governance requirements.

#### MECHANISM: GOVERNANCE REQUIREMENTS AS GRANT CONDITIONS

Funding agencies can condition grants on institutional establishment and maintenance of veto governance structures:

1. **Dual-Use Research Identification:** Grant conditions can require that institutions identify research falling within designated dual-use categories and subject such research to veto governance.
2. **Veto Body Establishment:** Grant conditions can require that institutions establish veto committees with specified composition, documented procedures, and protected decision-making authority.
3. **Review Requirements:** Grants can require that dual-use research be subject to veto body review before publication, licensing, or partnership formation.
4. **Funding Ineligibility for Violations:** Grant conditions can specify that institutions violating veto governance lose eligibility for future funding.
5. **Transparency and Reporting:** Grants can require that institutions maintain public records of veto decisions.

#### IMPLEMENTATION EXAMPLE

Suppose the National Science Foundation conditions grants on institutional governance:

756 *“Institutions must establish a Dual-Use Research Governance Committee. Com-*  
757 *mittee must include at least 2 faculty, 1 ethicist/policy expert, and 2 representa-*  
758 *tives from potentially affected communities. All research supported by this grant*  
759 *involving autonomous systems must be reviewed by the Committee before publi-*  
760 *cation, licensing, or commercial partnership. Institutions must maintain public*  
761 *records of committee decisions.”*

#### 762 763 ADVANTAGES OF FUNDING AGENCY IMPLEMENTATION

764 Funding agency requirements offer several advantages:

- 766 • **Leverage:** Funding agencies exercise significant leverage over institutions through grant  
767 awards.
- 768 • **Scale:** Conditions imposed by major funding agencies affect research across multiple in-  
769 stitutions.
- 770 • **Precedent:** Funding agencies have successfully imposed governance requirements in other  
771 domains (IRB requirements for human subjects research).
- 772 • **Flexibility:** Funding agencies can develop requirements tailored to different research do-  
773 mains.  
774

#### 775 776 LIMITATIONS AND CHALLENGES

777 Funding agency implementation faces challenges:

- 779 • **Political Opposition:** Requirements may face political opposition from researchers, uni-  
780 versities, and defense contractors.
- 781 • **International Coordination:** Without international coordination, national funding re-  
782 quirements become less effective.
- 783 • **Non-Compliance and Evasion:** Institutions might nominally comply while maintaining  
784 governance structures incapable of actually constraining militarization.
- 785 • **Institutional Autonomy:** Universities may resist funding conditions perceived as govern-  
786 mental overreach.  
787

### 788 789 C.4 REPOSITORY-LEVEL ARTIFACT GOVERNANCE: 790 INFRASTRUCTURE-BASED ENFORCEMENT

#### 791 792 CURRENT LANDSCAPE

793 AI research artifacts—code, models, datasets—are increasingly released through online repositories:  
794 Hugging Face, GitHub, ModelHub, arXiv, Zenodo. These repositories serve as primary distribution  
795 channels for research artifacts. Repositories currently function as neutral distribution channels with-  
796 out systematic governance of downstream uses.  
797

#### 798 799 MECHANISM: MACHINE-READABLE USE RESTRICTIONS AND LICENSING

800 Repositories can require that artifacts include machine-readable use restrictions specifying permitted  
801 and prohibited uses:

- 802 1. **License Standardization:** Artifacts can be released under licenses (based on GPL, MIT,  
803 or custom licenses) that explicitly specify use conditions. License text can be standardized  
804 using SPDX expressions.
- 805 2. **Metadata Embedding:** Models and datasets can include metadata fields specifying use  
806 conditions and governance requirements.  
807
- 808 3. **Access Restrictions:** Repositories can implement access control tiers. Some artifacts  
809 might be available for research access, while deployment access requires additional au-  
thorization.

- 810 4. **Version Control and Update Requirements:** Repositories can require that organizations  
811 deploying artifacts under specified conditions maintain current versions.  
812

813 MECHANISM: AUDIT AND DELISTING  
814

815 Repositories can implement compliance monitoring:  
816

- 817 1. **Deployment Reporting:** Users accessing artifacts for deployment can be required to sub-  
818 mit deployment context information.  
819  
820 2. **Audit Rights:** Repositories can assert rights to audit artifact uses, requiring deploying  
821 organizations to provide documentation.  
822  
823 3. **Violation Response:** If repositories determine that artifacts are being deployed in violation  
824 of use conditions, repositories can delist artifacts or revoke access.  
825  
826 4. **Public Transparency:** Repositories can maintain public records of use violations and en-  
827 forcement actions.

828 ADVANTAGES OF REPOSITORY-LEVEL GOVERNANCE  
829

830 Repository governance offers advantages:

- 831 • **Infrastructure Leverage:** Repositories already control access and distribution.  
832  
833 • **Technical Feasibility:** Machine-readable metadata and access control are established tech-  
834 nical capabilities.  
835  
836 • **Practical Friction:** Requiring users to certify compliance creates practical friction that  
837 increases consideration of restrictions.  
838  
839 • **Scalability:** A single repository governance decision affects all users accessing that repos-  
840 itory.  
841  
842 • **Transparency:** Repository-level governance is relatively transparent.

843 LIMITATIONS AND CHALLENGES  
844

845 Repository-level governance faces challenges:

- 846 • **Enforcement Limits at National Borders:** Repositories cannot enforce restrictions on  
847 users outside their legal jurisdiction.  
848  
849 • **Code Forking and Redistribution:** Repositories cannot prevent users from downloading  
850 artifacts, modifying them, and redistributing outside repository control.  
851  
852 • **Scientific Knowledge Cannot Be Restricted:** Repositories can restrict specific code or  
853 model implementations but cannot restrict underlying scientific knowledge.  
854  
855 • **Monitoring Burden:** Auditing deployments at scale requires significant resources.  
856  
857 • **User Deception:** Organizations can misrepresent intended uses or use restricted artifacts  
858 while evading deployment reporting.

859 C.5 COMMUNITY-BASED NORMS AND PROFESSIONAL SOLIDARITY  
860

861 CURRENT LANDSCAPE

862 Professional research communities have established ethics statements and codes of conduct. How-  
863 ever, these statements typically lack enforcement mechanisms and are developed within elite aca-  
ademic institutions without meaningful participation from affected communities.

864 MECHANISM: PROFESSIONAL COMMITMENTS AND COLLECTIVE ACCOUNTABILITY

865 Veto governance can be operationalized through professional norms and collective community com-  
866 mitments:

- 867
- 868 1. **Professional Codes of Conduct:** Research professional associations can establish bind-  
869 ing codes of conduct specifying that members commit to refusing military funding or to  
870 imposing use restrictions on dual-use work.
  - 871 2. **Institutional Commitments:** Research institutions can establish public commitments to  
872 governance of dual-use research, specifying that they will not accept certain types of mili-  
873 tary funding.
  - 874 3. **Community Oversight and Accountability:** Research communities can establish trans-  
875 parency mechanisms where members can raise concerns about other members' research or  
876 institutional practices.
  - 877 4. **Community-Led Governance Design:** Research communities should center affected  
878 communities in governance design, treating governance research as legitimate scholarly  
879 work.
  - 880 5. **Sanctuary Provisions:** Research communities can establish commitments to protect re-  
881 searchers who refuse military funding.
- 882

883

884 ADVANTAGES OF COMMUNITY-BASED GOVERNANCE

885 Community-based approaches offer advantages:

- 886
- 887 • **Legitimacy from Within:** Governance emerging from research communities has legiti-  
888 macy that external structures may lack.
  - 889 • **Flexibility and Adaptation:** Professional norms can be refined through community delib-  
890 eration.
  - 891 • **Low Institutional Burden:** Does not require new governmental authority or institutions.
  - 892 • **Centering Affected Communities:** Professional norms can explicitly center affected com-  
893 munities.
  - 894 • **Peer Accountability:** Professional peers understand research and can make nuanced judg-  
895 ments.
- 896

897

898 LIMITATIONS AND CHALLENGES

899 Community-based governance faces challenges:

- 900
- 901 • **No Enforcement Authority:** Absent institutional backing, professional norms are not  
902 binding.
  - 903 • **Freeriding and Competition:** Individual researchers or institutions that violate codes gain  
904 competitive advantages.
  - 905 • **Professional Heterogeneity:** Not all researchers identify with professional codes or asso-  
906 ciations.
  - 907 • **International Variance:** Professional norms vary across countries and research communi-  
908 ties.
  - 909 • **Limited Reach to Practitioners:** Professional codes primarily constrain academic re-  
910 search.
- 911

912

913 COMBINING IMPLEMENTATION PATHWAYS

914 The five implementation pathways are most effective when combined. A comprehensive approach  
915 might involve:

- 916
- 917 • Funding agencies condition grants on institutional governance

- Universities embed veto clauses in licenses and establish governance committees
- Conferences require dual-use disclosure
- Repositories implement use restriction metadata and compliance monitoring
- Professional communities establish codes of conduct and protect researchers exercising veto authority

## D DESIGN PRINCIPLES FOR LEGITIMATE VETO GOVERNANCE

Veto power can be abused. Governance structures designed without safeguards against abuse can become tools for suppression of legitimate research, for politicization of governance, or for discrimination against vulnerable populations. This section articulates design principles that constrain veto authority while preserving its capacity to prevent militarization.

### D.1 SCOPE LIMITATION AND PREDEFINED THRESHOLDS

#### PROBLEM: VETO WITHOUT BOUNDARIES

Unconstrained veto authority is dangerous. If veto bodies can refuse research for any reason or can continuously expand the scope of research subject to veto, veto becomes a mechanism for suppression rather than constraint on militarization. Historical examples demonstrate this risk: governments have used “national security” justifications to suppress inconvenient research, to silence critics, and to prevent investigation of governmental wrongdoing.

#### DESIGN PRINCIPLE: SCOPE LIMITATION

Legitimate veto governance must define veto authority narrowly. Veto should apply to specific downstream uses, not to research itself, publication, or inquiry:

1. **Downstream Use Limitation:** Veto authority should apply to transfer and deployment, not to research or publication. Researchers should retain freedom to conduct inquiries and to publish findings.
2. **Predefined Risk Categories:** Veto authority should apply only when research meets predefined risk criteria. Risk categories might include:
  - Integration into autonomous weapons lacking meaningful human control
  - Deployment in mass surveillance targeting political opponents or ethnic minorities
  - Use in violation of international humanitarian law or human rights law
  - Deployment in systems designed to enable forced displacement or ethnic cleansing
3. **High Confidence Thresholds:** Veto should be triggered by credible, documented risks, not speculative or distant possibilities.
4. **Exemptions for Legitimate Uses:** Veto should explicitly exempt legitimate civilian uses even if technologies are theoretically dual-use.

#### IMPLEMENTATION EXAMPLE

A veto body governing autonomous systems research might establish predefined thresholds:

“Veto Triggered If: Research has direct application to autonomous weapons lacking meaningful human control AND the research was funded by a defense agency OR the researcher has been approached by defense organizations about deployment OR the publication explicitly discusses weapons applications. Veto Not Triggered For: Research with general applicability to autonomous systems absent specific evidence of weapons application OR research with clear civilian applications absent evidence of military diversion.”

## D.2 COLLECTIVE DECISION-MAKING AND PROCEDURAL LEGITIMACY

### PROBLEM: VETO CONCENTRATED IN INDIVIDUAL AUTHORITY

Veto power concentrated in individual authority (a single researcher, administrator, or official) is dangerous. Individuals may make decisions based on personal bias, political motivation, or institutional self-interest rather than on principled assessment of risk.

### DESIGN PRINCIPLE: COLLECTIVE DECISION-MAKING

Legitimate veto governance should make veto decisions through deliberative collective processes:

1. **Pluralistic Veto Bodies:** Veto decisions should be made by committees including researchers with technical expertise, ethicists, legal experts, and representatives from affected communities. Pluralistic composition ensures that decisions reflect multiple perspectives.
2. **Documented Procedures:** Veto bodies should operate according to documented procedures specifying how veto is initiated, what deliberation occurs, what information is reviewed, how decisions are made, and what records are maintained.
3. **Transparency and Publicity:** Veto decisions should be documented and made public (with appropriate protections for sensitive information).
4. **Deliberation Requirements:** Veto bodies should be required to articulate reasoning for veto decisions.

### IMPLEMENTATION EXAMPLE

A university veto committee charter might specify:

*“Composition: 3 faculty researchers, 2 ethicists/policy experts, 2 community representatives, 1 legal expert. All members have equal voting authority. Decision Threshold: Veto requires 5 of 8 votes. Initiation: Veto can be initiated by any committee member, any faculty member, or external community representatives. Deliberation Timeline: 30 days from initiation. Documented Decision: Decisions recorded in writing including risk factors, harm prevention, and any conditions. Public Records: Decisions made public.”*

## D.3 FORMAL APPEAL AND REVISABILITY

### PROBLEM: VETO AS PERMANENT SUPPRESSION

Veto decisions that are permanent and unreviseable can become suppression mechanisms. If veto can never be overturned, veto transitions from governance to censorship.

### DESIGN PRINCIPLE: APPEAL AND REVISABILITY

Legitimate veto governance should make veto decisions subject to appeal and revision:

1. **Appeal Procedures:** Veto decisions should be subject to appeal to a higher authority independent of the original veto body.
2. **Grounds for Appeal:** Researchers should be able to appeal on grounds including procedural error, factual disagreement, applicability of predefined thresholds, or proportionality.
3. **Conditional Veto:** Veto decisions can be conditioned on circumstances under which veto could be lifted.
4. **Sunset Provisions:** Veto decisions can include sunset provisions specifying that the veto will be revisited at a future date.
5. **Evidence-Based Revision:** If new evidence becomes available, veto decisions should be revised based on evidence.

1026 IMPLEMENTATION EXAMPLE

1027

1028 A university might establish an appeals procedure:

1029

1030 *“First Appeal: Researchers can appeal veto committee decisions to a university*  
1031 *appeals board including external experts and community representatives. The ap-*  
1032 *peals board can reverse committee decisions if it finds procedural error or factual*  
1033 *error. Sunset: All veto decisions are revisited annually. If circumstances have*  
1034 *changed or evidence suggests veto is ineffective, the committee can revise deci-*  
1035 *sions.”*

1036

## 1037 D.4 LEGAL AND CAREER PROTECTIONS

1038

1039 PROBLEM: RETALIATION AGAINST THOSE EXERCISING VETO AUTHORITY

1040

1041 If researchers and administrators exercising veto authority face retaliation, veto authority remains  
1042 theoretical rather than practical.

1043

1044 DESIGN PRINCIPLE: PROTECTED REFUSAL

1045

1046 Legitimate veto governance must include legal and career protections:

1047

- 1048 1. **Anti-Retaliation Provisions:** Federal and institutional law should explicitly prohibit retal-  
1049 iation against researchers or administrators for exercising veto authority.
- 1050 2. **Tenure Protections:** Tenured faculty exercising veto authority should be protected by  
1051 tenure.
- 1052 3. **Whistleblower Protections:** Researchers or administrators reporting violations of veto  
1053 authority should be protected as whistleblowers.
- 1054 4. **Indemnification:** Institutions and individuals should be indemnified against legal liability  
1055 for good-faith exercise of veto authority.
- 1056 5. **Funding Protection:** Funding agreements should specify that institutional compliance  
1057 with veto governance does not affect grant funding.
- 1058 6. **International Protection:** For researchers in countries with weak legal protection, inter-  
1059 national professional organizations and funding agencies should provide legal support.

1060

## 1061 E ADDRESSING OBJECTIONS AND LIMITATIONS

1062

1063 Veto governance will face objections from multiple directions. This section addresses common  
1064 objections and acknowledges genuine limitations that veto governance does not resolve.

1065

### 1066 E.1 POLITICAL CAPTURE AND WEAPONIZATION

1067

1068 OBJECTION

1069

1070 “Veto power will be captured by hostile actors and weaponized against legitimate research. Govern-  
1071 ments will use veto mechanisms to suppress research investigating governmental wrongdoing.”

1072

1073 RESPONSE

1074

1075 This objection identifies a genuine risk. Governance mechanisms designed to constrain militariza-  
1076 tion could be misused to suppress inconvenient research. History provides examples: authoritarian  
1077 governments have used “national security” justifications to suppress research on environmental con-  
1078 tamination, governmental violence, and human rights abuses.

1079

1079 However, this risk must be weighed against the current absence of constraints on militarization.  
Currently, militarization proceeds largely unchecked. The choice is not between perfect governance

1080 and the status quo; the choice is between imperfect governance with some risk of abuse and no  
1081 governance at all with certainty of ongoing militarization.

1082 Several design features constrain risks of capture and abuse:  
1083

- 1084 • **Narrow Scope:** Veto authority limited to specific military uses
- 1085 • **Predefined Triggers:** Veto triggered by specific, documented risks
- 1086 • **Pluralistic Decision-Making:** Veto decisions made by committees including multiple per-  
1087 spectives
- 1088 • **Independent Review:** Veto decisions subject to appeal
- 1089 • **Transparency:** Veto decisions documented and made public
- 1090 • **International Observers:** International observers in governance processes
- 1091 • **Sunset Provisions:** Veto decisions time-limited and revisited

## 1092 E.2 INNOVATION STALLING AND RESEARCH SLOWING

### 1093 OBJECTION

1094 “Veto power will slow research and reduce innovation. Researchers will hesitate to pursue dual-use  
1095 work. The resulting slowdown will harm economic competitiveness and human welfare.”

### 1096 RESPONSE

1097 This objection reflects a particular value ranking: that speed of research and technological capability  
1098 are paramount. This is defensible but should be made explicit and should be subject to debate.

1099 Several points address this objection:

- 1100 • **Scope Limitation:** Veto applies only to downstream military use, not to research itself.  
1101 Researchers remain free to conduct inquiries and publish findings.
- 1102 • **Modest Time Impact:** Even if veto review requires time, the time impact should be modest  
1103 (30–60 days in research timelines spanning years).
- 1104 • **Selective Application:** Veto applies only to research with clear dual-use risks, not to main-  
1105 stream AI research.
- 1106 • **Competitive Fairness:** If veto governance is applied universally, all researchers face equiv-  
1107 alent constraints.
- 1108 • **Defensible Value Tradeoff:** There is defensibility to prioritizing safety and justice over  
1109 speed. The post-World War II international community decided that preventing nuclear  
1110 proliferation justified accepting slower development of nuclear energy technology.
- 1111 • **Harm Prevention Value:** Prevention of military AI integration creates real value: preven-  
1112 tion of autonomous weapons reduces civilian casualties, prevention of surveillance systems  
1113 targeting vulnerable populations protects autonomy and dignity.

## 1114 E.3 DEMOCRATIC LEGITIMACY OF VETO BODIES

### 1115 OBJECTION

1116 “Veto bodies will be undemocratic, wielding authority over research without legitimate democratic  
1117 sanction. Who appointed these veto committees?”

### 1118 RESPONSE

1119 This objection identifies a genuine tension: governance authority should be legitimate and account-  
1120 able. Several design features address democratic legitimacy concerns:

- 1134 1. **Community Leadership:** Veto bodies should be led and controlled by affected communi-  
1135 ties.
- 1136 2. **Elected or Appointed with Accountability:** Veto body members can be elected by re-  
1137 search communities or appointed through transparent processes with documented account-  
1138 ability.
- 1139 3. **Appeal and Revision:** Veto decisions should be subject to appeal and revision.
- 1140 4. **Public Participation:** Veto bodies should hold public meetings, accept public input, and  
1141 maintain public records.
- 1142 5. **Embedded in Democratic Institutions:** Veto governance should be embedded in demo-  
1143 cratic institutions (universities, funding agencies) that have governance mechanisms and  
1144 accountability structures.
- 1145 6. **Temporal Limits:** Veto body members should serve fixed terms with term limits.
- 1146
- 1147
- 1148

## 1149 F MARGINALIZED COMMUNITY PERSPECTIVES

1150 Veto governance succeeds as governance only if it is designed with and controlled by those most  
1151 harmed by military AI deployment. This section details why affected community leadership is  
1152 essential and how governance should be structured to center community voice and authority.

### 1153 F.1 EXCLUSION FROM CURRENT AI GOVERNANCE

#### 1154 THE GEOGRAPHY OF AI HARM

1155 AI governance discourse is concentrated in Global North institutions. Yet the most severe harms of  
1156 military AI deployment are borne in other regions. Autonomous weapons are tested and deployed in  
1157 conflicts in the Middle East, East Africa, and Central Asia. Mass surveillance systems are deployed  
1158 in countries with authoritarian governments. The populations experiencing these harms are largely  
1159 excluded from AI governance design.

1160 The exclusion of affected communities from governance design is both unjust and epistemically  
1161 limiting. Affected communities possess detailed, situated knowledge about:

- 1162 • How military AI systems are deployed and what harms result
- 1163 • What governance structures would be effective in their contexts
- 1164 • What unintended consequences governance mechanisms might have
- 1165 • How governance mechanisms intersect with other power dynamics affecting their commu-  
1166 nities

#### 1167 EPISTEMIC CONSEQUENCES OF EXCLUSION

1168 Researchers and policymakers designing governance from elite institutions lack situated knowledge.  
1169 Governance designed without affected community input risks being ineffective, culturally inappro-  
1170 priate, or inadvertently harmful.

### 1171 F.2 SHIFTING GOVERNANCE LEADERSHIP: OPERATIONALIZING COMMUNITY 1172 CONTROL

#### 1173 PRINCIPLE: AFFECTED COMMUNITIES LEAD

1174 Justice in governance requires that affected communities control governance design and implemen-  
1175 tation. This is not merely about inclusion or consultation; it is about control.

1188 OPERATIONALIZATION REQUIREMENTS

1189

1190 Shifting governance leadership to affected communities requires:

1191

1192 1. **Funding for Governance Research:** Governance research on military AI should be funded  
1193 and led by researchers in affected regions. Funding structures should be reoriented to sup-  
1194 port governance research in affected regions.

1195 2. **Treating Governance as Scholarship:** Governance research should be recognized as le-  
1196 gitimate academic work. Researchers conducting governance work should be eligible for  
1197 professorial positions and prestigious publication venues.

1198 3. **Binding Partnerships:** Governance partnerships should establish binding decision author-  
1199 ity for affected communities, not merely advisory roles.

1200 4. **Resource Control:** Affected communities should control governance resources—funding,  
1201 staffing, decision-making authority.

1202 5. **International Structures:** International governance structures should be led by represen-  
1203 tatives from affected regions with binding authority.

1204

1205 F.3 MILITARIZATION AND MARGINALIZATION: STRUCTURAL  
1206 INTERCONNECTIONS

1207

1208 DEEP LINKAGES

1209

1210 Militarization and marginalization are not separate problems; they are deeply interconnected. Mili-  
1211 tary violence disproportionately targets marginalized populations. AI systems amplify these dynam-  
1212 ics by automating and scaling targeting logics that already discriminate.

1213

1214 TARGETING DYNAMICS

1215

1216 Military violence targets populations perceived as threats. Populations perceived as threats are often  
1217 marginalized groups: ethnic minorities, indigenous peoples, religious minorities, political dissi-  
1218 dents. Military AI systems scale and automate these targeting dynamics.

1219 Addressing militarization without centering marginalization treats symptoms rather than structural  
1220 causes. Effective governance must address both military applications of AI and the way military AI  
1221 automation amplifies existing targeting and marginalization.

1222

1223 G CASE STUDIES OF MILITARIZATION PATHWAYS

1224

1225 This section details three historical cases where research transitioned into military applications,  
1226 identifying veto points where governance could have constrained militarization.

1227

1228 G.1 FACIAL RECOGNITION AND MASS SURVEILLANCE

1229

1230 RESEARCH DEVELOPMENT

1231

1232 Facial recognition technology developed over decades through open academic research. Major al-  
1233 gorithmic innovations came from researchers in academia and industry, with algorithms published  
1234 at conferences and in journals.

1235

1236 MILITARY PATHWAYS

1237

1238 Facial recognition technology moved into military and surveillance contexts through multiple path-  
1239 ways:

1240

1241 1. **Surveillance Deployment:** Law enforcement and military organizations began deploying  
facial recognition systems.

- 1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250
2. **Border Control:** Border control and immigration enforcement began deploying facial recognition systems.
  3. **Targeted Surveillance:** Authoritarian governments deployed facial recognition for targeted surveillance of ethnic minorities and political dissidents. Uyghur surveillance in Xinjiang relies heavily on facial recognition.
  4. **Weapons Integration:** Facial recognition systems began to be integrated into weapons systems for target identification.

#### 1251 VETO POINTS AND GOVERNANCE FAILURES

1252  
1253 Multiple veto points existed where governance could have constrained militarization:

- 1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264
1. **Publication Veto Point:** Individual papers could have been subject to veto review before publication or published with use restrictions.
  2. **Licensing Veto Point:** When systems were licensed to commercial entities or governments, licensing agreements could have included restrictions on surveillance use.
  3. **Partnership Veto Point:** University partnerships could have been subject to institutional review and could have been refused or conditioned.
  4. **Funding Veto Point:** Government funding could have been conditional on governance requirements.

#### 1265 WHY GOVERNANCE FAILED

1266  
1267 Governance failed at each veto point for structural reasons:

- 1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276
- No publication review of dual-use risks
  - No technology transfer governance requiring use restriction assessment
  - No partnership review mechanism
  - No funding conditions on governance
  - Competitive pressure to publish rapidly, license broadly, and accept military funding

## 1277 G.2 DARPA FUNDING IN AUTONOMOUS SYSTEMS

### 1278 RESEARCH DEVELOPMENT

1279  
1280  
1281 DARPA funds significant portions of AI research, particularly in autonomous systems, robotics, and sensor technology.

### 1282 MILITARY PATHWAYS

1283  
1284  
1285  
1286 DARPA funding explicitly aims to develop military capabilities. Much DARPA research transitions directly into military systems:

- 1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295
1. **Autonomous Vehicles:** Research funded by DARPA has transitioned into autonomous military ground and aerial vehicles.
  2. **Drone Control:** Research on drone control systems has transitioned into weapons systems.
  3. **Robotic Manipulation:** Research on robotic systems has transitioned into military robotics.
  4. **Autonomous Weapons:** Some DARPA funding explicitly aims to develop autonomous weapons capabilities.

1296 VETO POINTS AND GOVERNANCE FAILURES

1297

1298 Veto points existed in DARPA funding decisions:

1299

- 1300 1. **Funding Decision Veto Point:** Universities could have refused DARPA funding or condi-  
1301 tioned acceptance on governance requirements.
- 1302 2. **Institutional Governance Veto Point:** Universities could have established governance  
1303 committees reviewing DARPA-funded research before acceptance.
- 1304 3. **Publication Veto Point:** Universities could have restricted publication of DARPA-funded  
1305 research.
- 1306 4. **Technology Transfer Veto Point:** Technology could have been restricted from transfer to  
1307 weapons applications.

1308

1309 WHY GOVERNANCE FAILED

1310

1311 Governance failed for structural reasons:

1312

- 1313 • Revenue dependence on research funding
- 1314 • No institutional governance structures reviewing defense funding
- 1315 • Normalization of military research funding
- 1316 • Competitive pressure to accept DARPA funding
- 1317 • Lack of institutional authority to refuse military funding

1318

1319

1320 G.3 UNIVERSITY–DEFENSE CONTRACTOR PARTNERSHIPS

1321

1322 PARTNERSHIP DEVELOPMENT

1323

1324 Universities frequently establish research partnerships with defense contractors (Lockheed Martin,  
1325 Raytheon, Boeing, Northrop Grumman). These partnerships involve joint research, researcher visits,  
1326 technology transfer agreements, and collaborative publication.

1327

1328 MILITARY PATHWAYS

1329

1330 University-defense contractor partnerships transition research directly into weapons systems:

1331

- 1332 1. **Technology Transfer:** Technologies developed in partnerships are transferred to defense  
1333 contractors and integrated into weapons systems.
- 1334 2. **Researcher Expertise:** University researchers collaborate with defense contractors on  
1335 weapons development.
- 1336 3. **Credibility and Legitimacy:** Defense contractors benefit from association with university  
1337 researchers and academic legitimacy.
- 1338 4. **Talent Pipeline:** Partnerships create pipelines enabling university researchers to transition  
1339 into defense contractor positions.

1340

1341 VETO POINTS AND GOVERNANCE FAILURES

1342

1343 Veto points existed for partnership decisions:

1344

- 1345 1. **Partnership Approval Veto Point:** Universities could have established governance com-  
1346 mittees reviewing proposed partnerships.
- 1347 2. **Technology Transfer Veto Point:** Technologies could have been restricted from transfer  
1348 to weapons applications.
- 1349 3. **Researcher Participation Veto Point:** Researchers could have been prohibited from par-  
ticipating in weapons development partnerships.

1350 WHY GOVERNANCE FAILED

1351

1352 Governance failed for structural reasons:

1353

1354 • University partnerships provide funding and prestige

1355 • Researcher autonomy norms prioritize freedom to choose collaborators

1356 • No institutional governance structures reviewing partnerships

1357

1358 • Defense contractor participation in university governance

1359 • Cold War legacy normalizing university-defense partnerships

1360

1361 H OPERATIONALIZATION MECHANISMS AND ENFORCEMENT

1362 STRUCTURES

1363

1364 This section provides detailed mechanisms for operationalizing veto governance, addressing com-  
1365 mon concerns that veto governance lacks concrete implementation pathways.

1366

1367 H.1 ARTIFACT-LEVEL TRIGGERS AND DECISION PROTOCOLS

1368

1369 PROBLEM: VAGUE GOVERNANCE WITHOUT DECISION POINTS

1370

1371 Abstract veto authority is difficult to exercise. Without concrete decision points and specific proce-  
1372 dures, veto authority remains theoretical.

1373

1374 SOLUTION: PREDEFINED ARTIFACT-LEVEL TRIGGERS

1375

1376 Veto governance becomes practical when embedded in concrete decision points where choices must  
1377 be made.

1378

1379 H.0.1 GRANT ACCEPTANCE TRIGGER

1380

1381 **Decision Point:** When institutions receive funding from defense agencies (DARPA, military re-  
1382 search offices) or from other entities explicitly requesting dual-use research, a veto trigger activates.

1383 **Required Process:**

1384

1385 1. **Written Risk Memo** (Responsibility: Technology Transfer Office or Grants Administra-  
1386 tor)

1387 • Identify foreseeable military integration pathways

1388 • Assess likelihood and severity of militarization

1389 • Document evidence supporting risk assessment

1390 • Specify conditions that might mitigate risks

1391 • Timeline: 10 days from grant notification

1392 2. **Decision Deadline** (Responsibility: Veto Body)

1393

1394 • Veto body must convene and reach decision within 30 days

1395 • Failure to reach decision defaults to permission

1396 • Decision is recorded with reasoning

1397

1398 H.0.2 LICENSING TERMS TRIGGER

1399 **Decision Point:** When research is licensed to downstream organizations, a veto trigger activates.

1400

1401 **Required Process:**

1402

1403 1. **Written Risk Memo** (Responsibility: Technology Transfer Office)

• Identify proposed licensee and anticipated uses

- 1404
- Assess military applications likelihood
  - 1405 • Specify what use conditions might prevent militarization
  - 1406 • Timeline: 10 days from licensing proposal
  - 1407
- 1408 2. **License Agreement Terms** (Responsibility: Technology Transfer Office with Veto Body
- 1409 Input)
- License must include use restrictions approved by veto body
  - 1410 • License must require annual deployment reporting
  - 1411 • License must grant audit rights
  - 1412 • License must include indemnification clause for violations
  - 1413

1414 H.0.3 DEPLOYMENT CONTRACTS TRIGGER

1415

1416 **Decision Point:** When research transitions to deployment in actual use contexts, a veto trigger

1417 activates.

1418 **Required Process:**

1419

- 1420 1. **Deployment Impact Assessment** (Responsibility: Deploying Organization and Licensor)
- 1421
- Document deployment context
  - 1422 • Assess alignment with use restrictions
  - 1423 • Timeline: Assessment before deployment commences
  - 1424
- 1425 2. **Veto Body Review**
- Assess whether deployment violates licensing terms
  - 1426 • Approve, require modifications, or refuse deployment
  - 1427 • Timeline: 30 days from assessment submission
  - 1428

1429 H.0.4 REPOSITORY RELEASES TRIGGER

1430

1431 **Decision Point:** When research artifacts are released to public repositories, a veto trigger activates.

1432

1433 **Required Process:**

1434

- 1435 1. **Use Condition Specification** (Responsibility: Authors/Repository Maintainers)
- Artifacts released with machine-readable metadata
  - 1436 • License text specifies permitted and prohibited uses
  - 1437 • Repository implements access controls enforcing restrictions
  - 1438
- 1439 2. **Transparency and Monitoring** (Responsibility: Repository Platform)
- Repository maintains records of artifact access
  - 1440 • Repository flags high-risk artifact deployments
  - 1441 • Repository publishes aggregate compliance data
  - 1442
  - 1443

1444 H.2 VETO BODY COMPOSITION: STRUCTURED REPRESENTATION AND

1445 FAIRNESS

1446

1447 PROBLEM: ELITE CAPTURE AND MARGINALIZED EXCLUSION

1448

1449 Without structured composition requirements, veto bodies risk becoming dominated by privileged

1450 actors while excluding affected communities and marginalized voices.

1451

1452 SOLUTION: FIXED REPRESENTATION QUOTAS AND SUPPORT STRUCTURES

1453

1454 H.0.5 REPRESENTATION REQUIREMENTS

1455

1456 A representative veto body of 8–10 members should have mandatory composition (see Table 4):

1457

**Critical Principle:** All members have equal voting authority. Community representatives have full decision-making authority on par with faculty and experts.

Member Category	Minimum	Selection Method
Faculty Researchers (Technical)	2–3	Elected by research community
Ethicists / Policy Experts	1–2	Appointed by ethics/policy faculty
Affected Community Representatives	2–3	Appointed by affected communities
Legal Experts	1	Appointed by law faculty
External Experts (Optional)	1–2	Appointed by professional organizations

Table 4: Mandatory veto body composition

#### H.0.6 COMMUNITY REPRESENTATIVE SUPPORT STRUCTURES

Material support for community participation:

1. **Stipends:** Community representatives receive substantial stipends (e.g., \$50–100/hour or annual stipends of \$5,000–10,000).
2. **Access to Independent Legal Counsel:** Community representatives can access independent legal counsel paid by the institution.
3. **Travel and Logistics Support:** All costs covered; institutions cover all logistics for community members.
4. **Training and Capacity Building:** Institutions provide training helping community representatives understand technical issues and governance procedures.
5. **Accessibility Accommodations:** Flexible meeting times, remote participation options, interpreters for language/disability access.

#### H.0.7 TERM LIMITS AND ROTATION

1. **Term Length:** Veto body members serve fixed terms (2–3 years).
2. **Term Limits:** Members can serve maximum 2 consecutive terms (4–6 years total).
3. **Staggered Rotation:** New members appointed on staggered schedules.
4. **Community Control of Appointments:** Affected community organizations control appointment of community representatives.

### H.3 ENFORCEMENT MECHANISMS AND COMPLIANCE STRUCTURES

#### PROBLEM: VETO WITHOUT CONSEQUENCES

Veto authority without consequences is unenforceable. Enforcement requires that violations trigger automatic consequences.

#### SOLUTION: MULTI-LAYERED ENFORCEMENT ARCHITECTURE

#### H.0.8 CONTRACTUAL ENFORCEMENT

**Mechanism:** Research, licensing, and partnership agreements include explicit veto clauses creating legal obligations.

#### Example Veto Clause Language:

“Licensee agrees that:

- Research will not be deployed for autonomous weapons lacking meaningful human control; mass surveillance targeting political opponents; deployment in violation of international humanitarian law
- Licensee will not sublicense to parties likely to use for prohibited purposes
- Licensee will provide annual deployment context reports
- Licensee will permit licensor to audit deployment contexts

- 1512                   • Any violation entitles licensor to [remedies: license termination, indemnifi-  
1513                   cation for damages, injunctive relief]  
1514                   „

1515

#### 1516 H.0.9 MACHINE-READABLE USE RESTRICTIONS AND REPOSITORY ENFORCEMENT

1517  
1518 **Mechanism:** Artifacts include machine-readable use restriction metadata that repositories can auto-  
1519 matically enforce.

##### 1520 **Implementation:**

1521

1522 1. **License Metadata Standards:** Artifacts include metadata fields specifying use restrictions  
1523 using standard formats (SPDX expressions, YAML metadata).

##### 1524 2. **Repository Enforcement:**

- 1525                   • Flag artifacts with use restrictions  
1526                   • Require users to certify compliance before downloading  
1527                   • Track which organizations access restricted artifacts  
1528                   • Block access by organizations with known violations  
1529                   • Delist artifacts if violations detected

1530

#### 1531 H.0.10 COMPLIANCE AUDITING AND REPORTING

1532

1533 **Mechanism:** Organizations deploying research under use restrictions are subject to mandatory re-  
1534 porting and audit.

##### 1535 **Reporting Requirements:**

1536

##### 1537 1. **Annual Deployment Reports:**

- 1538                   • Deployment location and organization  
1539                   • Intended use and actual use  
1540                   • Populations affected by deployment  
1541                   • Measures ensuring compliance  
1542                   • Modifications or updates to deployed systems

##### 1543 2. **Audit Notification:**

- 1544                   • Facility visits to assess deployment context  
1545                   • Document review (deployment logs, system configurations)  
1546                   • Staff interviews about deployment practices  
1547                   • Written audit report

1548 3. **Public Reporting:** Licensor publishes aggregate data on compliance, deployment sectors,  
1549 violation rates, and enforcement actions.

1550

#### 1551 H.0.11 FUNDING AGENCY ENFORCEMENT

1552

1553 **Mechanism:** Federal funding agencies make compliance with veto governance a condition of grant  
1554 eligibility.

##### 1555 **Funding Conditions:**

1556

1557 1. **Institutional Governance Requirement:** Institutions must establish and maintain veto  
1558 governance structures. Failure triggers loss of funding eligibility.

1559 2. **Compliance with Veto Decisions:** Institutions must comply with veto decisions. Viola-  
1560 tions trigger investigation and potential defunding.

1561 3. **Transparency Requirement:** Institutions must maintain public records of veto decisions.

1562 4. **Loss of Funding for Non-Compliance:** Non-compliant institutions lose eligibility for fu-  
1563 ture grants.  
1564  
1565

1566 **Example Federal Regulation:**

1567

1568 “Award recipients must establish institutional governance procedures for dual-use  
1569 research. Research involving [specified dual-use categories] is subject to review  
1570 by an institutional committee meeting criteria specified in this regulation. Insti-  
1571 tutions must comply with committee recommendations regarding use restrictions,  
1572 licensing conditions, or research refusal. Institutions violating these requirements  
1573 become ineligible for future awards until compliance is restored.”

1574

1575 H.0.12 INDEMNIFICATION AND LEGAL LIABILITY

1576 **Mechanism:** Violations of use restrictions trigger financial and legal consequences.

1577

1578 **Indemnification Clause Example:**

1579

1580 “If Licensee deploys research in violation of use restrictions and that deployment  
1581 results in harm, Licensee indemnifies Licensor against:

- 1582 • Damages paid to harmed parties
- 1583 • Legal costs and attorney fees
- 1584 • Institutional liability
- 1585 • Regulatory fines or sanctions
- 1586 • Reputational harm

1587 ”

1588

1589 **Legal Liability:**

1590

- 1591 1. **Civil Liability:** Deploying organizations can be sued for damages resulting from viola-  
1592 tions.
- 1593 2. **Regulatory Liability:** Organizations face regulatory investigation and sanctions if deploy-  
1594 ment violates standards.
- 1595 3. **Criminal Liability:** Criminal prosecution may be possible in some jurisdictions if viola-  
1596 tions result in harm.

1597

1598 H.4 INTEGRATION WITH EXISTING INSTITUTIONAL INFRASTRUCTURE

1599

1600 KEY INSIGHT

1601

1602 The machinery for veto governance already exists within institutional structures. Implementation  
1603 requires embedding veto authority into existing decision-making processes.

1604

1605 H.0.13 UNIVERSITIES: REPURPOSING EXISTING AUTHORITY

1606

1607 **Technology Transfer Offices** already review licensing agreements and maintain licenses. **Mecha-**  
1608 **nism:** TTO charters updated to require:

- 1609 • Identification of dual-use research
- 1610 • Inclusion of use restrictions in licenses
- 1611 • Compliance monitoring and reporting
- 1612 • Authority to refuse licenses lacking adequate restrictions

1613

1614 **Faculty Governance Structures** already review research and establish institutional policies. **Mech-**  
1615 **anism:** Governance charters updated to:

1616

- 1617 • Require establishment of veto committee for dual-use research
- 1618 • Grant veto committee decision-making authority
- 1619 • Protect veto committee decisions from administrative override

1620 H.0.14 FUNDING AGENCIES: LEVERAGING GRANT CONDITIONS  
1621

1622 **Funding agencies** already establish grant conditions and monitor institutional compliance. **Mecha-**  
1623 **nism:** Grant conditions updated to:

- 1624 • Require institutional veto governance as condition of funding
- 1625 • Require compliance with veto governance decisions
- 1626 • Require transparency through public records
- 1627
- 1628

1629 H.0.15 REPOSITORIES: LEVERAGING PLATFORM INFRASTRUCTURE  
1630

1631 **Repositories** already manage metadata and control access. **Mechanism:** Repository platforms ex-  
1632 tend capabilities to:

- 1633 • Include use restrictions in standardized metadata
- 1634 • Implement access controls based on use restrictions
- 1635 • Automate enforcement of machine-readable restrictions
- 1636
- 1637

1638 H.0.16 PROFESSIONAL ASSOCIATIONS: ESTABLISHING NORMS  
1639

1640 **Professional associations** already establish codes of conduct and run conferences. **Mechanism:**  
1641 **Associations:**

- 1642 • Establish explicit commitments to veto governance in codes of conduct
- 1643 • Require dual-use disclosure for conference papers
- 1644 • Require use condition specification for published research
- 1645
- 1646

1647 NOTE  
1648

1649 Veto governance is not theoretical. The mechanisms detailed in this appendix show that veto gover-  
1650 nance can be operationalized through existing institutional structures, with concrete decision points,  
1651 structured procedures, material enforcement mechanisms, and appropriate safeguards against abuse.

1652 **Implementation requires political will, funding, and culture change.** But the mechanisms are  
1653 available. Universities can embed veto clauses in licenses. Funding agencies can condition grants  
1654 on governance. Repositories can implement use restrictions. Professional communities can establish  
1655 codes of conduct.

1656 *Affected communities must lead this transition.* Governance designed without affected community  
1657 leadership risks failing to address actual harms and risks replicating power imbalances that produced  
1658 militarization.

1659 The operationalization mechanisms detailed in this appendix provide the practical pathway for that  
1660 transition: from abstract commitment to concrete practice, from ethical exhortation to enforceable  
1661 governance, from elite-designed governance to community-led accountability.

1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673