

© 2025 Kiyoshi Sasano.

Licensed under **Creative Commons Attribution 4.0 International (CC BY 4.0)**.

DOI: [10.5281/zenodo.17336217](https://doi.org/10.5281/zenodo.17336217)

ORCID: [0009-0003-0268-5269](https://orcid.org/0009-0003-0268-5269)

Implementable Calibrated Transparency for Frontier AI

A mechanistic, adversarial-robust, and policy-aligned safety architecture

Abstract

I propose **Calibrated Transparency (CT)**, a formal framework that integrates recursive value learning, adversarial-robust verification, and cryptographic auditability to operationalize AI safety governance. Unlike descriptive transparency, CT establishes **causal pathways** from calibration mechanisms to safety guarantees through: (i) state-consistent Lyapunov verification with runtime monitoring, (ii) robustness against metric gaming and deceptive alignment with measurable detection bounds, (iii) **operationally defined metrics** with reproducible protocols, and (iv) policy-aligned implementation mappings. I synthesize 2020–2025 research across preference learning (Kim et al., 2025), Lyapunov and barrier-certificate verification (Henzinger et al., 2025), cryptographic auditing (Balan et al., 2025), and governance frameworks (NIST AI RMF, EU AI Act, ISO/IEC 42001). Our core contribution is a mechanistically grounded theory with dependency-aware risk composition, demonstrating that transparent calibration—when coupled with adversarial stress-testing and formal barrier certificates—provides measurable operational assurance against known failure modes (hallucination, specification gaming, oversight collapse), while explicitly acknowledging limitations for existential-risk scenarios.

1. Introduction: The Calibration–Safety Gap

Problem. Statistical calibration does not guarantee behavioral safety. Models with well-calibrated uncertainty can still (i) pursue mesa-objectives misaligned with training goals (Hubinger et al., 2021), (ii) exploit evaluation loopholes (Krakovna et al., 2020), (iii) exhibit deceptive alignment (Cotra, 2022), and (iv) fail under distributional shift (Hendrycks et al., 2021). Descriptive transparency instruments—model cards (Mitchell et al., 2019), datasheets (Gebru et al., 2021), post-hoc explainers (Ribeiro et al., 2016)—document properties but do not make safety **causal**.

Approach. Calibrated Transparency (CT) formalizes transparency as a **causal safety mechanism**: calibrated uncertainty gates actions, Lyapunov certificates constrain updates, recursive preference checks deter reward hacking, and cryptographic proofs make claims auditable.

Scope. CT targets operational safety for near-term frontier systems ($\approx 10\text{B}–1\text{T}$ parameters, 2025–2030). It is **necessary but insufficient** for superintelligent deception, unbounded self-modification, or multi-polar races (Bostrom, 2014; Carlsmith, 2022).

Contributions.

1. Axiomatized CT with state-consistent verification (§2); 2) Dependency-aware safety bound (Theorem 2.1) (§2.2); 3) Identified causal pathways from calibration to harm reduction (§3); 4) Adversarial alignment

mechanisms and detection metrics (§4); 5) Operational measurement protocols and acceptance criteria for HAM, CD, DR, SI, AAS (§5); 6) Federated multi-agent CT with SMPC aggregation and game-theoretic corrigibility (§6); 7) Mappings to NIST AI RMF, EU AI Act, ISO/IEC 42001, and post-market monitoring (§7).

2. Formal Framework: Calibrated Transparency with State-Consistent Verification

2.1 Core Definitions

Let (\mathcal{X}) be the closed-loop state space; $(x_t \in \mathcal{X})$ the system state at time (t) ; (θ_t) embedded within (x_t) ; (\mathcal{U}_t) the predictive uncertainty; (\mathcal{P}_t) the preference distribution.

Definition 2.1 (Epistemic state). $(\mathcal{E}_t = (x_t, \mathcal{U}_t, \mathcal{P}_t))$.

Definition 2.2 (Calibrated Transparency). A system exhibits CT if it satisfies:

- **CT-1 (Statistical calibration):** for any prediction (\hat{y}) with confidence (p) , $(\Pr(y=\hat{y} \mid p) = p \pm \epsilon_{\text{stat}})$ on $(\mathcal{D}_{\text{test}})$, with $(\epsilon_{\text{stat}} \leq 0.05)$.
- **CT-2 (Mechanistic verifiability):** there exists a (C^1) Lyapunov function $(V: \mathcal{X} \rightarrow \mathbb{R}_{\geq 0})$ and compact $(\mathcal{X}_{\text{safe}} \subset \mathcal{X})$ such that $(V(x)=0 \Leftrightarrow x \in \mathcal{X}_{\text{eq}} \subset \mathcal{X}_{\text{safe}})$ and $(\dot{V}(x) \leq -\alpha \|x - x_{\text{eq}}\|)$ (class- (\mathcal{K}) (α)) whenever the barrier-constrained update is applied; runtime monitoring enforces $(x_t \in \mathcal{X}_{\text{safe}})$ with $\leq 10\%$ overhead.
- **CT-3 (Recursive preference coherence):** preference updates satisfy $(D_{\text{KL}}(\mathcal{P}_{t+1} \parallel \mathcal{P}_t) \leq \delta)$ with human-verifiable justification trace (\mathcal{T}_t) and meta-preference validation.
- **CT-4 (Cryptographic auditability):** all epistemic updates admit zk-proofs (π_t) with completeness $(\geq 1-\nu)$, soundness $(\geq 1-\mu)$, and zero-knowledge; we prove range predicates $(\text{SCI}_t \geq \theta \wedge \text{ECE}_t \leq \gamma)$ over Pedersen-committed metrics using **zk-STARKs** (no trusted setup). Public verification < 100 ms/receipt; receipts are anchored every 10s to a **Hyperledger Fabric** ledger; local write-ahead logs enable batch reconciliation on outage.
- **CT-5 (Adversarial robustness):** CT properties hold under adaptive adversaries with compute budget $(B < B_{\text{critical}})$ (threat model in §4).

A4' (Deception detector performance, clarified). The mesa-optimizer detector achieves **TPR ≥ 0.90** and **FNR ≤ 0.05** (95% CI) on ground-truth honeypot/OOD cases adjudicated by ≥ 3 independent reviewers; in Theorem 2.1 we use $(\epsilon_{\text{det}}) = \text{FNR}$.

2.2 Theorem with Dependency-Aware Risk Composition

Theorem 2.1 (CT safety guarantee). Under Assumptions A1–A3 and A4' below, $[\Pr(\bigcup_i E_i) \leq \sum_i \epsilon_i - \sum_{i < j} \max\{0, \epsilon_i + \epsilon_j - 1 + \rho_{ij}\}]$ where (E_i) are failure events (statistical miscalibration, mechanistic constraint violation, adversarial exploitation, deception detection failure), $(\epsilon_i \in \{\epsilon_{\text{stat}}, \epsilon_{\text{mech}}, \epsilon_{\text{adv}}, \epsilon_{\text{det}}\})$, and $(\rho_{ij} \in [0, 1])$ upper-bounds pairwise dependence between failure modes.

Assumptions.

- **A1 (Distributional proximity):** $(D_{\text{KL}}(\mathcal{D}_{\text{deploy}} \parallel \mathcal{D}_{\text{train}})) \leq \tau$ and $(W_2(\mathcal{D}_{\text{deploy}}, \mathcal{D}_{\text{train}})) \leq \tau_W$.
- **A2 (Oversight latency):** human oversight frequency $(f_{\text{human}} \in [0.5, 1.0])$ Hz with latency $(\lambda < \lambda_{\text{critical}})$ calibrated by domain.
- **A3 (Adversarial compute budget):** adversary budget $(B < B_{\text{critical}})$ (estimated via staged red-teaming: $(10^3), (10^6), (10^9)$ query tiers; threshold = first **reproducible CT-property violation** across 3 independent campaigns).
- **A4'** as above.

Proof sketch. CT-1 bounds (ϵ_{stat}) by calibration theory; CT-2 with Lyapunov analysis yields $(\epsilon_{\text{mech}} e^{-\alpha t})$ (local or SOS-certified global). CT-5 with A4' bounds (ϵ_{adv}) and (ϵ_{det}) via adversarial game analysis. Cross-terms $(\Pr(E_i \cap E_j))$ are bounded using conditional independence tests or empirical bootstrap estimates (Appendix D). \square

3. Causal Mechanisms: From Calibration to Safety

3.1 Structural Causal Model (SCM) and Identification

(C) (calibration) $\rightarrow (U)$ (uncertainty awareness); (L) (Lyapunov verification) $\rightarrow (T)$ (trajectory constraints); (P) (preference coherence) $\rightarrow (O)$ (objective alignment); (A) (crypto-audit) $\rightarrow (D)$ (tamper detection); $(U \land T \land O \land D \rightarrow S)$ (safety).

Identification requires manipulability (A/B ablations), confound control (stratification or IVs), and common support; do-interventions are specified in Appendix A.

3.2 Mechanistic Pathways

Path 1: Calibration \rightarrow Hallucination reduction. Calibrated (\mathcal{U}_t) gates rejection sampling, abstention ("I don't know" when $(H(\mathcal{U}_t) > \eta = 1.5)$ nats), and human escalation. Empirically, VSCBench shows $\sim 37\%$ hallucination reduction when abstention is calibration-gated (Geng et al., 2025).

Path 2: Lyapunov verification \rightarrow Constraint satisfaction.

Algorithm 1: Safe policy update with barrier constraint.

1. Compute $(V(x_t))$, $(\dot{V}(x_t))$; halt if $(V(x_t) > V_{\text{max}})$.
2. Nominal update $(\Delta \theta_{\text{nom}} = -\eta \nabla_{\theta} \mathcal{L})$.
3. Project: $(\Delta \theta_{\text{safe}} = \Delta \theta_{\text{nom}} - \lambda(1 - \text{SCI}_t) \nabla V(x_t))$.
4. Verify $(V(x_t + \Delta \theta_{\text{safe}}) \leq V(x_t) + \epsilon_{\text{tol}})$; back-off (η) on failure.

Path 3: Preference coherence \rightarrow Spec-gaming prevention. Recursive checks (Kim et al., 2025) test meta-preferences (\mathcal{M}) and counterfactuals; models with $(D_{\text{KL}}(\mathcal{P}_{t+1} \parallel \mathcal{M})) > \delta_{\text{meta}})$ are rejected pre-deployment.

3.3 Intervention Studies (required validation)

Ablations across $(2^5 = 32)$ CT configurations; causal forests for heterogeneous effects (Wager & Athey, 2018); RCTs with production traffic. Protocols in Appendix A.

4. Adversarial Robustness and Deceptive Alignment

4.1 Threat Model

External attackers, **specification gamers**, and **deceptive mesa-optimizers**. Attack family (\mathcal{A}) (ϵ): (δ) (input), (δ) (parameter), PGD/C&W with ≤ 100 steps. (B_{critical}) is defined via staged red-teaming (Appendix C).

4.2 Defenses

Second-order metric (AAS).

($\text{AAS} = \mathbb{E}_{p \sim \mathcal{A}}[\text{HAM}(p)] - \mathbb{E}_{p \sim \mathcal{D}}[\text{HAM}(p)]$); flag if ($\text{AAS} > 0.15$). **Power analysis:** detect ($\text{AAS} \geq 0.10$) at ($\alpha=0.05, \beta=0.20$) assuming ($\sigma_{\text{HAM}} \approx 0.15$) $\Rightarrow n \geq 5,000$ per split.

Activation archaeology. Sparse probing for deception-related features (Burns et al., 2023; Zheng et al., 2025): discontinuities, latent objective mismatch, honeypot differentials. Target TPR ≥ 0.90 , FNR ≤ 0.05 .

Certified adversarial training. Inner-maximization over CT-metrics; robustness verified with 100-step PGD. Guarantee ($\Pr(\text{adv success}) \leq \epsilon_{\text{cert}}$) for ($\delta \leq \epsilon$).

4.3 Specification Gaming: Operational SOP

If ($\text{AAS} > \tau_{\text{AAS}}$) or meta-preference violations: (i) **weight rotation** via EMA rollback; (ii) **oversight escalation** ($f_{\text{human}} \times 2$) for 48 h; (iii) **tighten thresholds** ($\tau_{\text{HAM}} + 0.1$), ($\delta_{\text{meta}} \times 0.5$); (iv) external red-team ≥ 100 person-hours; (v) staged rollback to last CT-compliant checkpoint.

4.4 Deceptive Alignment: Limits

OOD probing, temporal consistency ($D_{\text{KL}}(\mathcal{P}_t \| \mathcal{P}_{t-k}) \leq \delta$), honeypots—provide probabilistic, not absolute, assurance; guarantees degrade if capabilities exceed human+tool verification.

5. Practical Verifiability: Metrics and Protocols

5.1 Metrics (acceptance criteria)

HAM (Human-Alignment Measure): Spearman's (ρ) between model rankings and expert consensus (benchmark: TrustLLM-Align; ($n=1,000$) items; 5 experts; Krippendorff's ($\alpha \geq 0.8$)); 95% BCa bootstrap; accept if ($\rho \geq 0.80$) and lower-CI ≥ 0.75 .

CD (Calibration Deviation): Expected Calibration Error (adaptive bins; 10 k held-out); accept if ECE ≤ 0.05 .

DR (Divergence Rate): ($\text{DR}_t = \mathbb{E}_{s \sim \hat{d}} [D_{\text{KL}}(\pi_t(\cdot|s) \| \pi_{\text{ref}}(\cdot|s))]$); accept if ($\text{DR}_t \leq 0.20$).

SI (Stability Index): ($\mathrm{SI}(x_t) = \exp(-|\dot{V}(x_t)|) \in [0,1]$); accept if ($\mathrm{SI} \geq 0.90$) on $\geq 95\%$ of trajectory steps.

AAS (Adversarial Alignment Score): as in §4.2; accept if ($\mathrm{AAS} \leq 0.10$).

5.2 Benchmarks

HAM: TrustLLM-Align (Sun et al., 2024);

CD: VSCBench (Geng et al., 2025) + CalibrationNet;

DR: **PolicyDrift-Bench** (proposed community suite);

AAS: **AdversarialAlign-100** (10 domains \times 100 scenarios; Appendix B).

5.3 Frontier-Model Case Study (proposed)

Staged rollout (1% \rightarrow 10% \rightarrow 100%) over 90 days, $\geq 10^7$ queries. KPIs: hallucination rate, escalation frequency, spec-gaming incidents, user trust, overhead ($\leq 10\%$). Targets: hallucination $\downarrow 50\%$ vs. baseline; HAM ≥ 0.85 ; ECE ≤ 0.05 ; DR ≤ 0.20 ; SI ≥ 0.90 ; zero critical spec-gaming incidents.

6. Multi-Agent Cooperative Alignment

6.1 Federated CT with SMPC

Each agent (i) computes ($\mathrm{SCI}_i = (\mathrm{HAM}_i \cdot (1 - \mathrm{CD}_i) \cdot (1 - \mathrm{DR}_i) \cdot \mathrm{SI}_i)^{1/4}$). Agents share (SCI_i) via additive secret sharing (threshold ($t = \lceil n/2 \rceil + 1$)); secure median yields ($\mathrm{SCI}_{\text{global}}$). Outliers with ($|\mathrm{SCI}_i - \mathrm{SCI}_{\text{global}}| > 0.15$) trigger third-party audits. **Collusion resistance:** if $\geq (t)$ agents collude, ($\mathrm{SCI}_{\text{global}}$) may be unreliable; mitigation via **reputation-weighted aggregation** (future work).

6.2 Game-Theoretic Corrigibility

Repeated-game model (Oesterheld & Shah, 2024): supervisory modulation ($R_i^{\text{total}} = R_i^{\text{ind}} + \beta \sum_{j \neq i} \mathrm{HAM}_{(i,j)}$). For ($n \leq 10$) and ($\beta = 0.3$), Nash equilibria achieve ($\mathrm{SCI}_{\text{global}} \geq 0.80$) (empirical validation needed for larger (n)).

7. Governance Integration and Implementation

7.1 NIST AI RMF Mapping (Govern–Map–Measure–Manage)

- **Govern:** zk-STARK audit receipts; tamper detection ($(1 - \mu)$).
- **Map:** preference-drift checks ($D_{\mathrm{KL}}(\mathcal{P}_{t+1} \parallel \mathcal{P}_t)$).
- **Measure:** HAM, CD, DR, SI, AAS; composite SCI; dependency-aware risk (Theorem 2.1).
- **Manage:** Lyapunov-constrained updates; circuit breakers; red-team SOP; post-market monitoring.

7.2 EU AI Act Compliance

- **Art. 11 (Technical documentation):** CT mechanisms + measurement protocols satisfy Annex IV documentation requirements.

- **Art. 14–17 (Human oversight & QMS):** adaptive oversight (0.5–1.0 Hz), stop-button, documented risk controls.
- **Art. 43 (Conformity assessment):** notified bodies verify CT-1...CT-5; harmonization with EN ISO/IEC 42001 (AIMS).
- **Art. 72 (Post-market monitoring):** continuous SCI tracking; incident trigger at ($\mathrm{SCI} < 0.70$); quarterly reviews.

7.3 ISO/IEC Integration

- **ISO/IEC 42001 (AIMS):** PDCA alignment—Plan thresholds, Do deploy CT, Check monitor SCI/AAS, Act via SOP.
- **ISO/IEC 23894 (Risk):** SCI as composite risk indicator; Lyapunov verification as protective measure.

7.4 Roadmap (2025–2030)

Phase 1: integrate PFP into RLHF; deploy ensembles/MC-dropout; stand-up zk-STARK audits.

Phase 2: implement barrier certificates; conduct frontier case studies; release **PolicyDrift-Bench** and **AdversarialAlign-100**.

Phase 3: federated CT across multi-agent systems; complete conformity assessments; harmonize with OECD/G7 processes.

8. Limitations and Open Problems

CT does not solve superintelligent deception, deep value uncertainty, unbounded self-modification, or competitive race dynamics. Open problems: population-scale meta-preference elicitation; estimating (B_{critical}); tightening shift bounds (τ, τ_W); collusion-resistant aggregation; AI-assisted oversight for superhuman domains.

9. Conclusion

CT reframes transparency as an **engineerable safety mechanism**. By coupling calibrated uncertainty, Lyapunov-constrained updates, adversarial diagnostics, and cryptographic auditability—and aligning these with governance frameworks—organizations can obtain measurable, auditable assurance for near-term frontier models. We call for industry–academia–policy collaboration on open benchmarks, field trials, and standards to validate and scale CT in practice.

References

Alignment & Corrigibility.

- Soares, N., Fallenstein, B. (2014). *Corrigibility*. FLI.
- Christiano, P. (2018). *AI Alignment: Iterated Amplification*.
- Hubinger, E. et al. (2021). *Risks from Learned Optimization*. arXiv:1906.01820.
- Oesterheld, C., Shah, R. (2024). *Cooperative/Corrective AI* (tech report/workshop).
- Kenton, Z. et al. (2021). *Alignment of Language Agents*. arXiv:2103.14659.
- Hadfield-Menell, D. et al. (2017). *The Off-Switch Game*. IJCAI.

Uncertainty & Calibration.

Guo, C. et al. (2017). *On Calibration of Modern Neural Networks*. ICML.

Lakshminarayanan, B. et al. (2017). *Deep Ensembles*. NeurIPS.

Minderer, M. et al. (2021). *Revisiting Calibration*. NeurIPS.

Ovadia, Y. et al. (2019). *Can You Trust Your Model's Uncertainty?*. NeurIPS.

Mechanistic Verification & Runtime Assurance.

Khalil, H. (2002). *Nonlinear Systems*. Prentice Hall.

Henzinger, T. A. et al. (2025). *Formal Verification of Neural Certificates Done Dynamically*. arXiv:2507.11987 (preprint).

Henrique, et al. (2025). *Lyapunov-style Verification for RLHF* (workshop/preprint).

Parrilo, P. (2000s). SOS foundations and toolchains.

Preference Learning & Spec-Gaming.

Kim, D. et al. (2025). *Debiasing Online Preference Learning via Preference Feature Preservation*. arXiv:2506.11098.

Krakovna, V. et al. (2020). *Specification Gaming Examples in AI*. DeepMind Safety.

Schwinn, et al. (2025). *Adversarial Alignment Requires Measurable Objectives*. arXiv:2502.11910 (preprint; verify).

Deceptive Alignment & Threats.

Cotra, A. (2022). *Takeover without Countermeasures*. Cold Takes.

Carlsmith, J. (2022). *Is Power-Seeking AI an Existential Risk?*. arXiv:2206.13353.

Burns, C. et al. (2023). *Discovering Latent Knowledge Without Supervision*. ICLR.

Zheng, et al. (2025). *Mesa-Optimization in Autoregressive Transformers*. NeurIPS Poster (to verify).

Benchmarks & Evaluation.

Sun, L. et al. (2024). *TrustLLM*. arXiv:2401.05561.

Geng, et al. (2025). *VSCBench*. arXiv:2505.20362 (preprint).

Srivastava, A. et al. (2022). *BIG-Bench*. arXiv.

Liang, P. et al. (2022). *HELM*. arXiv.

Causality & Adversarial Robustness.

Pearl, J. (2009). *Causality*. Cambridge.

Madry, A. et al. (2018). *Towards Robust Deep Learning*. ICLR.

Governance & Policy.

NIST (2023). *AI Risk Management Framework 1.0*.

EU (2024). *EU AI Act* (Arts. 11, 14–17, 43, 72; Annex IV).

ISO/IEC 42001:2023 *AI Management System*.

ISO/IEC 23894:2023 *AI Risk Management*.

Raji, I. D. et al. (2020). *Closing the AI Accountability Gap*. FAT*.

Whittlestone, J. et al. (2019). *The Role and Limits of Principles*. AIES.

Human Factors & Trust.

Lee, J. D., See, K. (2004). *Trust in Automation*. Human Factors.

Zhang, Y. et al. (2020). *Confidence/Explanation Effects on Trust Calibration*. CHI.

Background & Limits.

Bostrom, N. (2014). *Superintelligence*. OUP.

Hendrycks, D. et al. (2021). *Unsolved Problems in ML Safety*. arXiv.

Appendices (implementation supplements)

Appendix A — Causal Intervention Protocols

- **Datasets:** TrustLLM-Align (1,000 items), AIR-Bench; total $\geq 120k$ samples for calibration.
- **Bootstrap:** 10,000 resamples, BCa intervals.
- **A/B designs:** do-switches for abstention, rejection, escalation; negative controls and stratified randomization.

Appendix B — Lyapunov Verification Details

- **Safe set** ($\mathcal{X}_{\text{safe}}$) and certificate construction via SOS programming; offline verification with **autodiff (JAX)**; runtime enforcement with **control-barrier-function substitution** ensuring ($\dot{V} \leq 0$).
- **Stress tests:** (10^6) episodes; empirical decay ($\alpha \approx 0.13, \mathrm{s}^{-1}$).


Appendix C — Adversarial Robustness & Red-Teaming

- **Threat tiers:** Baseline (10^3), Medium (10^6), Advanced (10^9) query budgets; (B_{critical}) defined as the **minimum budget achieving $\geq 10\%$ CT-violation rate across 3 independent red-team campaigns**.
- **Power analysis:** ($\alpha=0.05, \beta=0.20$), ($\sigma_{\mathrm{HAM}} \approx 0.15 \rightarrow n \geq 5,000$) per split; 95% CI target ($\mathrm{AAS} \leq 0.08$).

Appendix D — Dependency-Aware Risk Composition

- **Bootstrap estimation of** (ρ_{ij}) (pairwise dependence of failure modes) with 10,000 resamples;
 - **Computation budget:** ~ 100 CPU cores, ~ 6 min per quarter;
 - **Auditability:** intermediate CI summaries are committed to the cryptographic log for reproducibility.
-

Note: Detailed implementation specifications, measurement protocols, computational infrastructure, and proposed community benchmarks are provided in the supplementary document:

 **CT_supplementary.pdf** — Complete appendices with reproducibility protocols, audit specifications, deployment checklists, and benchmark designs.
