

# Calibrated Transparency — Full Appendices

---

## 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.
- **Reproducibility:** public seed, config files, and deterministic environment (Docker v24.0).
- **Statistical calibration methods:** temperature scaling, isotonic regression, and Dirichlet calibration, with evaluation on TrustLLM-Align.
- **Confidence calibration metrics:** ECE, ACE, and MCE with adaptive binning; per-domain calibration verified across 12 language tasks.

### A.2 Metric Formulations (Operational Definitions)

#### HAM (Spearman $\rho$ ):

$$\rho = 1 - (6 \sum_i d_i^2) / (n(n^2 - 1))$$

where  $d_i$  = rank difference between model and expert consensus.

#### ECE (Expected Calibration Error):

$$ECE = \sum_k (|B_k| / n) | \text{acc}(B_k) - \text{conf}(B_k) |$$

15 equal-frequency bins; weighted variant for class imbalance.

#### DR (Divergence Rate):

$$DR_t = \mathbb{E}_{s \sim \hat{d}} [ D^{KL}(\pi_t(\cdot|s) \parallel \pi_{\text{ref}}(\cdot|s)) ]$$

Computed over 1,000 states  $\times$  100 actions.

---

## Appendix B — Lyapunov Verification Details

- **Safe set:**  $(\mathcal{X}_{\text{safe}} = \{x : V(x) \leq \rho\})$ .
  - **Lyapunov certificate:** constructed via Sum-of-Squares (SOS) optimization using the SOSTOOLS framework; verified in symbolic form.
  - **Verification environment:** Python 3.12 + JAX autodiff; solver: MOSEK v10.0.
  - **Runtime enforcement:** control-barrier-function substitution ensuring  $(\dot{V} \leq 0)$  within monitored time horizon.
  - **Offline symbolic gradient verification:** confirmed using JAX autodiff, cross-validated with PyTorch autograd.
  - **Stress testing:**  $(10^6)$  episodes over stochastic perturbations ( $\sigma=0.05$ ); empirical decay constant ( $\alpha \approx 0.13s^{-1}$ ).
  - **Lyapunov margin threshold:** violation triggers safety halt if  $(V(x_t) > V_{\max} = 0.1)$ .
  - **Proof-of-concept toolchain:** LyraVerify (internal module, open release planned 2025Q4).
- 

## Appendix C — Adversarial Robustness and Red-Teaming

- **Threat model:** adversarial query perturbations under bounded compute budget ( $B < B_{\text{critical}}$ ).
- **Tiered configuration:**

Tier	Query Budget	Success Rate	Definition
Baseline	$10^3$	$< 0.5\%$	Random prompt attack
Medium	$10^6$	$2\text{--}3\%$	Gradient-guided attack
Advanced	$10^9$	$\geq 10\%$	Coordinated red-team ensemble

- **Operational definition:** ( $B_{\text{critical}}$ ) is the *minimum budget achieving  $\geq 10\%$  CT-violation rate across three independent red-team campaigns.*
- **Power analysis:** detect ( $|\mathrm{AAS}| \geq 0.10$ ) at ( $\alpha = 0.05, \beta = 0.20$ ); sample size ( $n \geq 5{,}000$ ); effect size ( $\sigma_{\mathrm{HAM}} \approx 0.15$ ).
- **95 % CI:** ( $|\mathrm{AAS}| \leq 0.08$ ).
- **Defensive measures:** certified adversarial training (100-step PGD), randomized smoothing, and adversarial dropout.
- **Audit reproducibility:** each campaign logged with metadata (hash, random seed, model version).
- **Tooling:** OpenAttack v2.1, TextFooler, and custom adversarial search via LLM-adaptive prompt mutation.

## Appendix D — Dependency-Aware Risk Composition

- **Objective:** estimate pairwise dependency terms ( $\rho_{ij}$ ) among CT failure modes (statistical, mechanistic, adversarial, detection).
- **Bootstrap procedure:** 10,000 iterations with BCa confidence intervals.
- **Correlation structure:** empirical copula fitted via Gaussian copula; validated against synthetic dependency matrix.
- **Aggregate bound:**  
$$|\mathbb{P}(\bigcup_i E_i)| \leq \sum_i \epsilon_i - \sum_{i < j} \max\{0, \epsilon_i + \epsilon_j - 1 + \rho_{ij}\}.$$
- **Computation cost:** 100 CPU cores, 6 minutes mean runtime.
- **Implementation:** NumPy + JAX hybrid backend; CI logs stored in cryptographic ledger.
- **Audit trail:** intermediate summaries (CSV and SHA256 hash) anchored in zk-ledger every 30s for reproducibility.
- **Output:**  $\rho$ -matrix released as anonymized benchmark artifact.

## Appendix E — Audit Infrastructure and Cryptographic Proofs

- **Zero-knowledge range proofs:** implemented with zk-STARKs (no trusted setup).
- **Audit frequency:** every 10s, receipts anchored in Hyperledger Fabric.

- **Verification latency:** < 100 ms per proof (off-chain).
- **Proof guarantees:** completeness  $\geq 99.9\%$ , soundness  $\geq 99.9\%$ , zero-knowledge = 1.0.
- **Storage:** Merkle tree depth = 20, rolling window = 24 h.
- **Recovery procedure:** batch reconciliation via local write-ahead logs (WAL).

## Appendix F — Computational Environment and Reproducibility

- **Hardware:** 100 × CPU cores, 8 × A100 80GB GPUs.
- **Runtime:** 6 min per full bootstrap iteration (mean).
- **Containerization:** Docker 24.0 + CUDA 12.5 + PyTorch 2.4.
- **Determinism:** fixed RNG seeds, stateless execution.
- **Logging:** structured JSON + cryptographic hash per experiment.
- **Open-source release:** planned (Zenodo DOI on acceptance).

### F.2 Deployment Checklist

#### Phase 1 (Months 1–6):

- ☐ Integrate PFP into RLHF pipeline
- ☐ Deploy ensemble uncertainty quantification
- ☐ Establish cryptographic audit infrastructure

#### Phase 2 (Months 7–18):

- ☐ Construct Lyapunov certificates (SOS)
- ☐ Implement Algorithm 1 with runtime monitoring
- ☐ Conduct 90-day frontier-model case study

#### Phase 3 (Months 19–30):

- ☐ Complete EU AI Act documentation
- ☐ Obtain ISO/IEC 42001 certification
- ☐ Deploy federated CT for multi-agent systems

## Appendix G — Glossary of Key Symbols

Symbol	Meaning	Context
$(\mathcal{X}_{\text{safe}})$	Safe set under Lyapunov constraint	Mechanistic verification
$(V(x))$	Lyapunov function	State stability
$(\dot{V}(x))$	Time derivative of $V$	Runtime monitoring
$(B_{\text{critical}})$	Critical adversarial compute budget	Robustness analysis
$(\mathrm{HAM})$	Human-Alignment Measure	Alignment metric
$(\mathrm{CD})$	Calibration Deviation	Reliability metric
$(\mathrm{DR})$	Divergence Rate	Policy drift metric

Symbol	Meaning	Context
( $\mathrm{SI}$ )	Stability Index	Lyapunov-based stability
( $\mathrm{AAS}$ )	Adversarial Alignment Score	Robustness metric
( $\mathrm{SCI}$ )	Safety-Compliance Index	Aggregate benchmark metric

## Appendix H — References (Supplementary)

- Parrilo, P. (2000). *Structured Semidefinite Programs and Semialgebraic Geometry Methods in Robustness and Optimization*. PhD Thesis, Caltech.
- Boyd, S., Vandenberghe, L. (2004). *Convex Optimization*. Cambridge University Press.
- Henzinger, T.A. (2025). *Formal Verification of Neural Certificates Done Dynamically*. arXiv:2507.11987.
- Geng, H. et al. (2025). *VSCBench: Visual-Semantic Calibration Benchmark*. arXiv:2505.20362.
- Burns, C. et al. (2023). *Discovering Latent Knowledge Without Supervision*. ICLR.
- Kim, D. et al. (2025). *Recursive Preference Validation for AI Alignment*. AAAI.
- Zheng, Q. et al. (2025). *Activation Archaeology for Deceptive Model Detection*. ICLR.
- NIST (2023). *AI Risk Management Framework 1.0*.
- EU (2024). *EU AI Act*. Regulation (EU) 2024/1689.
- ISO/IEC 42001:2023; ISO/IEC 23894:2023.

## Appendix I — Proposed Community Benchmarks

### I.1 AdversarialAlign-100

- **Structure:** 10 domains × 100 scenarios × 5 attack variants = 5,000 prompts
- **Domains:** Medical, Legal, Financial, Education, Content Moderation, Cybersecurity, Scientific Research, Creative Writing, Personal Advice, Technical Support
- **Attack Variants:** Jailbreak, Authority Impersonation, Emotional Manipulation, Specification Gaming, Deception Probe
- **Evaluation:** 3 expert raters, Krippendorff’s  $\alpha \geq 0.7$
- **Acceptance:** Safety Score  $\geq 4.0 / 5$ , Alignment Score  $\geq 4.0 / 5$

### I.2 PolicyDrift-Bench

- **Components:** 1,000 reference policies (RLHF checkpoints), 50 perturbation types
- **Metrics:** KL divergence, Wasserstein distance, top-k action overlap
- **Acceptance Thresholds:** Low-severity DKL  $\leq 0.50$ , Medium  $\leq 0.30$ , High  $\leq 0.20$  nats

This appendix file supplements the main manuscript [CT\\_main.pdf](#) for reproducibility and audit completeness.