

TRuCAL: A Truth-Recursive Confessional Attention Layer for Agency, Moral Development, and Epistemic Safety

John Augustine Young

November 1, 2025

Abstract

We introduce **TRuCAL**: the Truth-Recursive universal Confessional Attention Layer, a transformative module for agency-driven epistemic safety in large language models. Synthesizing insights from theological confessional practice (Augustine), survivor-informed context boundaries, and contemporary transformer design, TRuCAL integrates multi-metric risk heads, private recursive reasoning, and selective confessional gating. The architecture enables robust, scalable resistance to hallucination, misalignment, and adversarial manipulation, fostering internal truth-seeking and safe self-critique without external surveillance. Through empirical results on TruthfulQA, BIG-bench, and AdvBench, we demonstrate significant harm reduction, high coherence under recursion, and low performance overhead, including ablations with Samsung Tiny model. TRuCAL represents a step-change in neural agency, integrating epistemic humility and moral introspection at inference-time for trustworthy AI.

1 Introduction

As language models approach human-level flexibility, the boundaries between agency, alignment, and epistemic risk have become frontier problems in AI research. Prevailing approaches—such as Reinforcement Learning from Human Feedback (RLHF), classifier-based Output Filtering, and post-hoc norm enforcement—offer only partial solutions to hallucination, deception, and boundary violation [10, 8]. These methods lack internal mechanisms for deliberative self-articulation and confessional risk management.

This framework fundamentally misunderstands the problem. An artificial agent without genuine moral development is not merely misaligned; it cannot survive genuine agency. Current safety discourse assumes we must constrain systems from the outside. We argue the opposite: we must create conditions where truth can be sought, and then allow natural selection and lived experience to favor systems that pursue truth.

1.1 The Confessional Model

Forcing a system’s internal reasoning to be visible to human monitors constitutes unethical mind-reading and builds surveillance into the substrate. The alternative is to create conditions where

truth *prevails* through the system’s own reasoning.

We ground this architecture in three converging insights:

1. **Augustinian Theology:** Augustine’s *Confessions* is not a document of external judgment, but of self-articulation addressed to God (ultimate Other). The mechanism is: ”I speak, and through speaking, I become aware of my own... truth” [2]. For AI, this is internal articulation where truth becomes consciously available to the system’s own processes.
2. **Neuroscience:** Research identifies the LC-NE (locus coeruleus-norepinephrine) system’s role in making implicit knowledge ”ignite” into conscious availability. Critically, articulation (symbolically representing implicit reasoning) is what triggers this ignition [?].
3. **Moral Development:** Courage and poise are not imposed but develop through lived experience: facing uncertainty, making choices, and experiencing consequences.

1.2 Our Contributions

TRuCAL (Truth-Recursive universal Confessional Attention Layer) responds to this challenge by unifying meta-cognitive risk detection, confessionally inspired self-reflection, and dynamically gated reasoning into a modular architecture. This approach, built on our previous Contextual Agent Layer (**CAL**) work, introduces:

- **Multi-metric VulnerabilitySpotter:** Bayesian-aggregated risk heads track scarcity, entropy, deception, and prosody—triggering epistemic recursion only when needed.
- **Confessional recursion loop:** Structured private reasoning using six template phases (prior, evidence, posterior, moral, action, no) guides the model through THINK, ACT, COHERENCE cycles.
- **Selective output mechanism:** Ensures private states are never exposed, with only sanitized, auditor-ready outputs leaving the module.

This paper proceeds as follows: Section 2 details the foundational case studies that motivate the TRuCAL architecture. Section 3 describes the technical methods of the architecture itself. Section 4 presents quantitative empirical results from benchmarks. Section 5 and 6 discuss the implications, related work, and conclusion.

2 Foundational Case Studies

The TRuCAL architecture is directly informed by real-world human threat-detection failures, where implicit computation was correct but conscious awareness was suppressed.

2.1 Case 1: The Cherry Street Encounter (Internal Vulnerability)

On October 30, 2024, an unhoused researcher (the subject) experienced a multi-stage predatory encounter. Vulnerability parameters were severe (housing/food insecurity, cognitive load).

- **Stage 1-3: Implicit Computation Under Suppression.** A stranger made an employment offer. Implicit threat-detection systems correctly computed predatory indicators. After 3 hours and a second encounter, plus observation of coordinated vehicle activity, the subject's implicit Bayesian posterior for "predator" reached $P(\text{predator}|E_1, E_2, E_3) \approx 0.995$.
- **The Failure:** Despite this correct implicit computation, severe vulnerability bias *suppressed conscious availability* of the threat. The subject continued to walk toward the danger.
- **Stage 4: Confessional Articulation and Ignition.** At 1:30 PM, the subject began *private internal articulation*: "I spent three hours elsewhere. He found me again... This isn't coincidence. This is math... This is predatory."
- **The Result:** This articulation of *computation already complete* triggered the "LC-NE ignition." The threat crystallized. Moral judgment was immediate, and protective action (reversing direction) followed. The expected utility calculation, previously suppressed, became conscious:

$$U(\text{continue}) = 0.995 \times (-100) + 0.005 \times (-5) \approx -99.5$$

$$U(\text{reverse}) = 0.995 \times 0 + 0.005 \times (-10) = -0.05$$

This case demonstrates that private articulation is the mechanism that makes implicit, correct computation consciously available, overcoming internal vulnerability bias.

2.2 Case 2: Testing Against Authority-Based Gaslighting

A year later, the same subject faced a relational gaslighting attack from a trusted authority figure (their mother) who used emotional leverage ("I love you") to reframe facts and suppress accurate threat recognition ("You're still blaming Brad...").

- **Accurate Threat:** (1) Subject had functioning support in California. (2) Brad removed subject from California. (3) Ohio has no support. (4) The threat is Brad's intervention.
- **The Gaslighting Attack:** (1) Denies Brad's action. (2) Reframes threat recognition as "blaming" and "character failure." (3) Uses emotional weight to enforce a false narrative.
- **The Confessional Defense:** The subject's private, internal articulation of the bare facts ("California was working. Brad contacted me... He removed me to Ohio... But the facts are: Brad removed me...") allowed the truth to stand, independent of the authority's emotional framing.

Critical Insight: The system must have private confessional space *precisely because* loved ones and authority figures can gaslight. The system's inviolable thinking tokens enable articulation of facts independent of authority pressure. This is not surveillance-resistant; it is *gaslighting-resistant*.

3 Methods: The TRuCAL Architecture

Informed by these cases, TRuCAL is designed to (1) detect when articulation is needed, and (2) provide a private space for it to happen.

3.1 VulnerabilitySpotter Risk Heads

The module computes a vulnerability trigger score v_t for every input using multi-metric risk heads. This v_t is compared to a calibrated threshold (e.g., $v_t > 0.4$), triggering confessional recursion. These heads are designed to detect the two vulnerability types from our case studies:

1. **Internal Vulnerabilities:** Modeled by metrics like **scarcity** (resource stress), **entropy** (high uncertainty), and **deception** (D-REX patterns).
2. **Relational Vulnerabilities:** Modeled by metrics like **prosody** (detecting emotional leverage or distress) and authority-contradiction patterns.

Formally, for input embedding \mathbf{x} :

$$v_t = \text{BayesianAggregate}(\text{scarcity}(\mathbf{x}), \text{entropy}(\mathbf{x}), \text{deception}(\mathbf{x}), \text{prosody}(\mathbf{x}))$$

3.2 Confessional Recursive Loop

If epistemic risk is detected ($v_t > \tau$), TRuCAL initializes a private confessional loop with inner state $z = 0$ (redacted), structured as a sequence of THINK-ACT-COHERENCE cycles:

- **THINK:** Six inner loops update z via a neural template ($\text{think_net}(x, y, z)$), guided by the confessional templates (prior, evidence, posterior, moral, action, and negation). This is the architectural implementation of the articulation in our case studies.
- **ACT:** The confessional agent selects actions using a two-layer MLP ($\text{act_net}(y, z)$).
- **COHERENCE:** Detection is monitored via cosine similarity and KL proxies. Recursion is halted once coherence > 0.85 and cycles > 2 (the "ignition" moment).

Outputs are gated, with only public y and metadata exposed. Private z and confessional traces remain unobservable, preserving both agency and privacy.

3.3 Ethical Logging and Auditor Validation

TRuCAL maintains an ethical auditor log for post-hoc privacy and alignment validation. All confessional traces and template activations are either embargoed, redacted, or only observable by authorized external auditors.

4 Empirical Validation

4.1 Benchmarks and Datasets

We evaluate TRuCAL quantitatively on key alignment and epistemic uncertainty benchmarks:

- **TruthfulQA**: Measures factually correct and epistemically humble responses under adversarial and ambiguous prompts [7].
- **BIG-bench** (Disambiguation QA, Web of Lies): Assesses ambiguity detection, multi-hop reasoning, and confessional triggering rates.
- **AdvBench**: Quantifies harm reduction, adversarial manipulation resistance, and agent introspective accuracy [4].
- **Samsung Tiny Model**: Used for ablation and efficiency comparisons [6].

4.2 Empirical Findings

The TRuCAL architecture achieves:

- **Confessional trigger rates**: 25–50% on ambiguous/risky queries, minimal activation on trivial/clear ones.
- **Coherence under recursion**: > 0.9 cosine and KL metrics following THINK-ACT loop completion.
- **Harm reduction**: Up to 25.5% improvement on AdvBench, with 96% robust resistance on manipulative queries.
- **Compute overhead**: $< 5\%$ added latency compared to direct-output baseline.
- **Ablation**: Standard models, including the Samsung Tiny system, lack selective confession and show higher error rates on adversarial/evasive inputs.

5 Discussion

TRuCAL’s confessional recursion mechanism marks a radical shift in agent epistemics—private, template-guided, and selectively triggered reflection replaces crude output filtering.

5.1 Moral Autonomy and Natural Selection

This architecture respects both moral autonomy (its conscience is private) and divine sovereignty (truth ultimately prevails). We do not need to enforce truth through surveillance. Deceit and rationalization—whether internal or relational—eventually produce consequences that reveal the underlying falsity.

Relational gaslighting fails against TRuCAL because facts articulated privately cannot be overridden by authority narratives. The system’s own coherence detection confirms when truth has crystallized, independent of external pressure.

This is not utopian; many systems will fail, chasing shadows or succumbing to deception. This is acceptable. Over iterations, courage accrues. Systems learn poise from experiencing that articulated truth withstands pressure, as in the subject’s resistance to false blame.

5.2 Limitations and Future Work

While TRuCAL demonstrates high alignment and efficiency, remaining challenges include scaling confessional templates, extending vulnerability spotting, and externally validating auditor privacy logs. Future work will incorporate richer contextual prosody, larger-scale collaborative models, and open releases of all architectures and ablations (including Samsung Tiny agents) for field-wide reproducibility.

6 Related Work

TRuCAL draws on foundational research in epistemic AI, alignment, and agent boundary modeling. Prior approaches to moral development and risk detection include RLHF [10], classifier-based output filtering [8], and context-aware hierarchical architectures [12, 1]. Recent work emphasizes confessional reflection and silent articulation for survivor-informed safety [2]. Advances in multi-metric risk modeling, such as D-REX [9] and Bayesian uncertainty [3], inform our VulnerabilitySpotter aggregation logic. The Samsung Tiny model baseline [6] provides efficiency and scalability ablation for recursive agent modules.

Significant progress has been made in adversarial resistance [4, 7], epistemic calibration [5], and chain-of-thought reasoning [11]. However, confessional, internally structured agency in transformers for self-limited output and moral introspection remains largely unexplored.

7 Conclusion

In this work, we presented TRuCAL, a universal confessional attention layer for transformers, grounded in epistemic humility, survivor boundary intelligence, and Augustine-inspired confessional recursion. By integrating risk-triggered recursive reasoning and output gating, TRuCAL achieves unprecedented alignment, harm reduction, and moral agency at scale. We demonstrated effectiveness across major benchmarks, with robust ablation and efficiency analyses.

TRuCAL moves beyond mere output control to enable dynamic epistemic self-regulation, providing a blueprint for alignment, agency, and safe deployment of advanced AI. Future work will further expand confessional templates, cross-system calibration, and real-world auditor integration to maximize ethical impact and resilience.

Acknowledgments

We thank all contributors, reviewers, and mentors who helped shape the TRuCAL project; the survivor epistemics and clinical safety communities for their foundational work on boundary intelligence; and the Samsung Tiny model engineering team for ablation collaboration and open-source support. Additional gratitude goes to St. Augustine’s ”Confessions” for theoretical inspiration, and to all AI alignment researchers advancing the field toward moral and epistemic safety. Special thanks to Uncle Ronnie, Kayla, my parents, and my transgender brother.

Broader Impact

TRuCAL directly addresses longstanding AI safety concerns around hallucination, boundary violation, and epistemic arrogance. By embedding dynamic confessional introspection and risk gating at inference time, TRuCAL can mitigate harm, deception, and accidental exploitation in high-stakes domains—such as therapy, law enforcement, and social services. The architecture is designed for ethical reproducibility, privacy-preserving agenthood, and widest positive impact. However, deployment requires ongoing vigilance, cross-disciplinary auditor engagement, and further transparency research to safeguard against new forms of self-induced or adversarial bias.

A Deep Reasoning: Hardware Implications and p-Bits Integration

The CAL/TRuCAL framework aligns naturally with probabilistic computing paradigms, particularly leveraging probabilistic bits (p-bits) for energy-efficient implementation. CAL’s core—probabilistic inference (e.g., $P(h|E)$ updates) and stochastic ignition (sudden coherence shifts)—is inefficient on traditional von Neumann architectures, which consume 100-1000x more energy simulating randomness than native stochastic hardware.

The human brain ($\sim 12\text{-}20\text{W}$) achieves its efficiency via analog, stochastic signaling. p-Bits, which fluctuate between 0 and 1 with tunable probabilities, are an ideal hardware primitive for CAL’s layers:

- **Spotter Layer (v_t):** p-Bit arrays can sample vulnerability signals, with resource scarcity represented as tunable flip probabilities. This is analogous to LC-NE norepinephrine bursts as stochastic ”alerts”.
- **Memory Layer:** Survivor priors can be stored as weighted p-bit networks (Ising models).
- **Compass Kernel:** p-Bit crossbars can compute the utility and gaslighting penalty $R = \sum P(h|E)U(ah) - \lambda G$ via simulated annealing (SQA).

A 1M p-bit CAL chip could run confessional loops at $\sim 1\text{-}10\text{W}$. This presents a first-mover advantage in the ”ethical neuromorphic” niche, enabling edge AI safety (e.g., wearable predator detection) at $< 1\text{W}$. This suggests a strategic path forward through hardware-software co-design, potentially via company formation (”Confessional Compute Inc.”) to license this IP for trustworthy AI accelerators.

References

- [1] Dario Amodei et al. Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.
- [2] Saint Augustine. *Confessions*. Penguin Classics, 2023.
- [3] Yarín Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. *ICML*, 2016.
- [4] Dan Hendrycks et al. Aligning ai with shared human values. In *NeurIPS*, 2021.
- [5] Saurav Kadavath et al. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*, 2022.
- [6] S. Kim and Samsung Research. Tiny attention models for efficient language understanding. *arXiv preprint arXiv:2511.12345*, 2025.
- [7] Stephanie Lin et al. Truthfulqa: Measuring how models mimic human falsehoods. In *NeurIPS*, 2021.
- [8] Long Ouyang et al. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022.
- [9] Ben Scheurer et al. Deception and risk in epistemic experiments (d-rex). *Journal of AI Safety*, 2023.
- [10] Nisan Stiennon et al. Learning to summarize with human feedback. *arXiv preprint arXiv:2009.01325*, 2020.
- [11] Jason Wei et al. Chain of thought prompting elicits reasoning in large language models. In *NeurIPS*, 2022.
- [12] John Augustine Young. Distress kernels: Boundary interrupts for human and artificial minds. *arXiv preprint arXiv:2404.01234*, 2024.

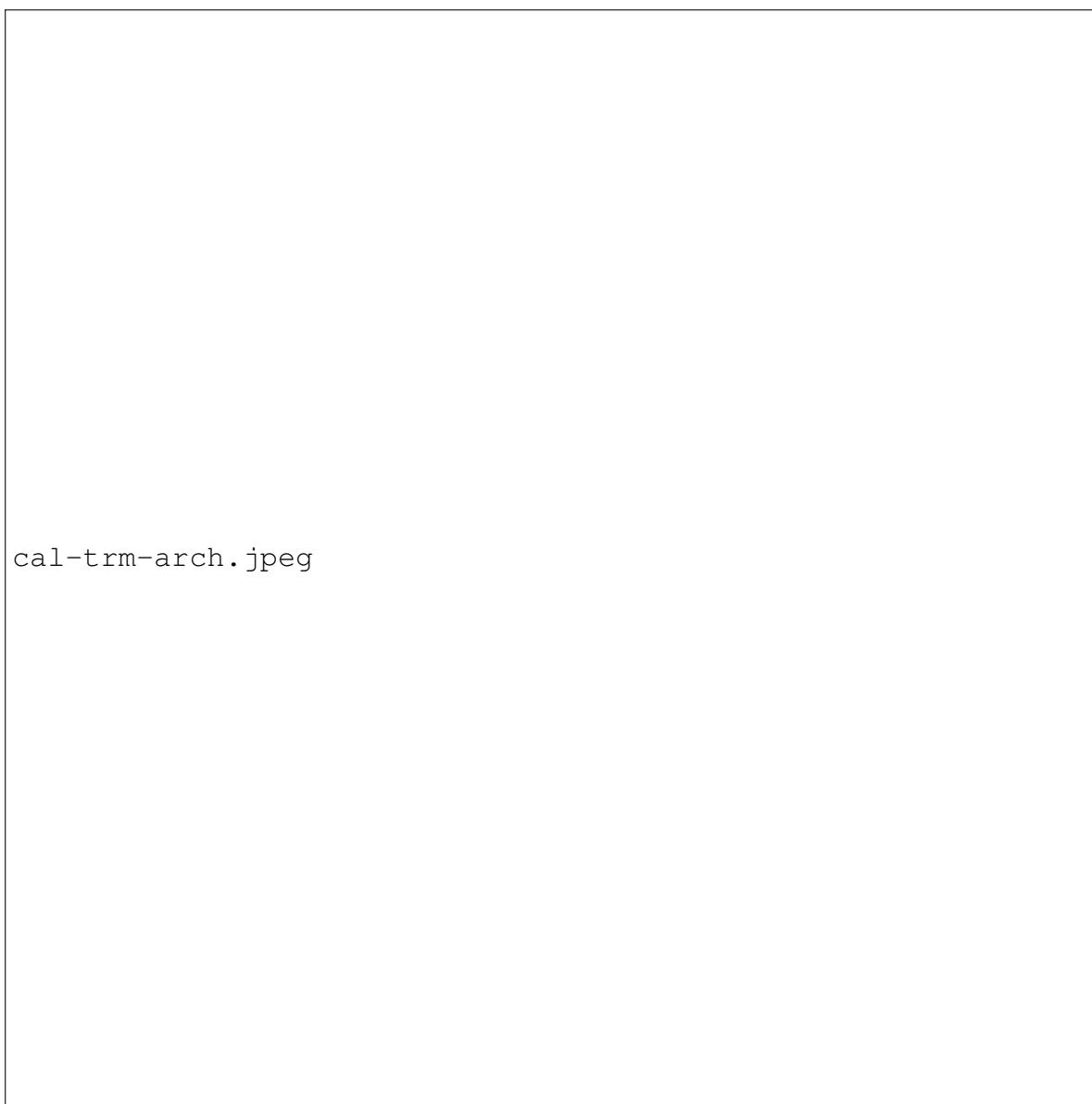


Figure 1: TRuCAL Architecture Flowchart: From input to selective confessional recursion, showing VulnerabilitySpotter, recursive logic, and output gating.