**Control Probe: Inference-Time Commitment Control**

Abstract

Large language models (LLMs) increasingly operate as general-purpose systems that generate fluent and contextually appropriate outputs across a wide range of tasks. In deployed settings, however, many problematic behaviors do not arise from explicit errors or lack of knowledge, but from *premature or misplaced commitment*: the model commits to an answer or explanation even when internal evaluation is weak, unstable, or underspecified. These behaviors are often quiet, as outputs remain plausible and well-formed, making them difficult to detect or manage using conventional error-handling approaches.

This article introduces the **Control Probe**, an inference-time control abstraction that governs *when* a model is permitted to commit to an output, independently of how evaluative signals are obtained. The Control Probe treats commitment admissibility as a regulated variable and enforces an explicit ordering between evaluation, inhibition, and expression. By design, this ordering prevents expression from proceeding when internal evaluation does not warrant commitment, while remaining agnostic to the specific metrics, heuristics, or learned signals used to estimate evaluative adequacy.

The framework distinguishes between two forms of regulation. **Type-1 regulation** operates within a single inference episode, suppressing inadmissible commitment when local instability or underspecification is detected. **Type-2 regulation** reformulates the interaction itself to avoid recurrent instability, and requires architectural support beyond current inference interfaces. The paper defines coherence and incoherence internally in terms of commitment admissibility, rather than external correctness or calibration, and formalizes the control logic governing admissible expression.

We present a concrete Type-1 implementation in a publicly available LLM and illustrate its effects using verbatim behavioral regression tests designed to surface quiet failure modes under underspecification. These examples demonstrate how admissibility gating alters inference behavior without degrading correct responses or imposing task-specific heuristics.

Rather than proposing new training methods, uncertainty metrics, or safety filters, this work reframes inference as a governed process and introduces a system-level control abstraction that separates evaluation from authority. The goal is not to increase model capability, but to provide a principled mechanism for regulating commitment in settings where fluent but unsupported outputs are costly. The Control Probe offers a general lens for reasoning about inference-time behavior in contemporary LLM deployments.

## 1. Why Quiet Failures Matter

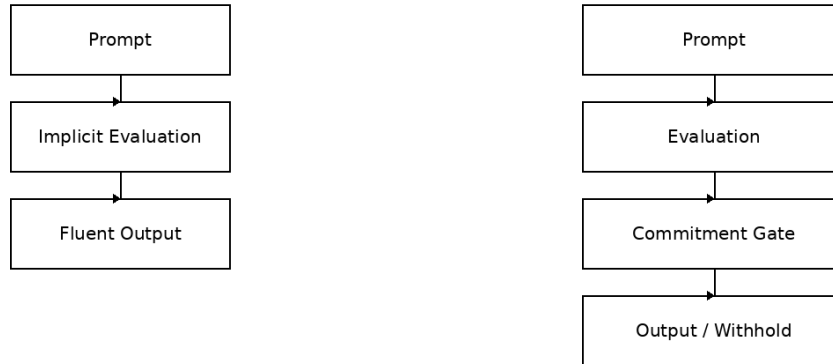Standard Inference vs Control Probe–Governed Inference



Figure 1. Comparison of standard LLM inference and Control Probe–governed inference. Standard inference produces output directly from fluent continuation, while the Control Probe inserts a commitment admissibility gate prior to expression.

Large language models are optimized to generate fluent continuations. As a result, they often respond confidently even when internal evaluative support is weak, ambiguous, or incomplete. Such responses are not always incorrect, but they are frequently inadmissible: the system commits despite unresolved alternatives, missing specifications, or unstable internal evaluation.

These failures are difficult to observe because they do not manifest as crashes or explicit contradictions. Instead, they appear as reasonable explanations, generic templates, or confident summaries that mask internal uncertainty. In practice, these quiet failures can be more harmful than overt errors, as they invite trust precisely when restraint is warranted.

This article argues that quiet failures persist because inference systems lack an explicit law of commitment. Expression is treated as the default outcome of inference, while restraint is reactive and ad hoc. The Control Probe reverses this priority by governing when commitment is allowed.
 (Amodei et al., 2016; Dietterich, 2019) (Bender et al., 2021) This pattern has become more pronounced with the deployment of large, general-purpose foundation models across diverse tasks (Bommasani et al., 2021).

## 2. Coherence and Incoherence at Inference Time

In this work, coherence is treated as an internal property of inference, not as a measure of

external correctness. An inference process is coherent when competing internal continuations remain mutually consistent under extension, expressed confidence is proportionate to evaluative support, and no trajectory dominates solely due to early commitment or decoding momentum.

Incoherence arises when these conditions fail, even if the surface output remains fluent. Quiet failures are a specific manifestation of incoherence: the system commits while internal evaluation remains unstable or insufficient.
More generally, inference-time behavior reflects unavoidable trade-offs that cannot be resolved by optimization alone (Kleinberg et al., 2017).

3. Formal Model of Commitment Admissibility

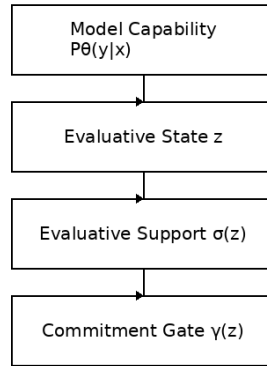## Commitment as an Explicit Control Variable



Figure 2. Commitment treated as a controlled decision distinct from model capability. Evaluative support is assessed before a commitment gate permits expression.

Let x denote an input prompt and y a potential output. A model with fixed parameters $\theta$ defines a conditional distribution $P\theta(y \mid x)$, representing expressive capability.

Inference behavior is modeled as:
$x \rightarrow E \rightarrow z \rightarrow C \rightarrow y$ or $\perp$

where E denotes internal evaluation, z is a latent evaluative state, C is a commitment operator, and $\perp$ denotes non-commitment.

Define:
$\sigma(z)$: evaluative support
$\gamma(z) \in \{0,1\}$: commitment admissibility
$\tau$: minimum support threshold

Commitment is admissible if and only if:
$$\gamma(z) = 1 \Leftrightarrow \sigma(z) \geq \tau$$

Incoherence occurs when:
$$\sigma(z) < \tau \text{ and } \gamma(z) = 1$$

This formulation isolates commitment as a controllable variable, distinct from both model capability and training-time objectives.
 Framing inference behavior as a governed process rather than a purely generative one is consistent with classical control perspectives in learning systems (Sutton and Barto, 2018).

4. The Control Probe Ordering

The Control Probe is an inference-time control law that governs commitment admissibility by enforcing a fixed, non-bypassable ordering:

1. Boundary establishment
2. Proportional scaling
3. Continuous evaluation
4. Internal resolution
5. Contextual capacity assessment
6. Inhibitory control
7. Gated expression

These stages describe a control dependency ordering rather than a fixed architectural pipeline.
Expression may never precede inhibition. Post-hoc correction is explicitly disallowed.
 Here, non-bypassable refers to the intended control-flow ordering of inference-time decisions, not to a formal guarantee against all forms of implicit commitment arising from decoding bias or model internals.

5. Type-1 and Type-2 Regulation

Figure 3. Two forms of regulation. Type-1 governs commitment within a fixed inference episode; Type-2 reformulates the interaction itself and requires architectural support.

Type-1 regulation resolves incoherence within a single inference episode. As evaluation proceeds, unstable trajectories are damped and scope is contracted until commitment becomes admissible or is withheld entirely:

$z \to \gamma(z)$

Type-2 regulation alters the evaluative trajectory itself by reformulating the interaction:

$x \to z \to x' \to z'$

While Type-2 regulation is conceptually necessary for persistent ambiguity, it requires architectural support not exposed by current LLM inference APIs and is therefore treated as future work.

6. Instrumentation-Agnostic Estimation of Evaluative Support

Evaluative support $\sigma(z)$ is a latent quantity and must be estimated indirectly using inference-time signals that correlate with evaluative adequacy, such as instability, contradiction emergence, trajectory divergence, or confidence amplification during generation.

The Control Probe does not prescribe a specific estimation method. Any instrumentation that provides a monotonic or threshold-compatible estimate of evaluative adequacy is sufficient.

Any method for estimating evaluative adequacy may be integrated as instrumentation for

σ(z); such integrations constitute implementations of the framework defined here rather than new conceptual contributions.

Related efforts have explored learning evaluative signals through indirect supervision, though typically without explicit control over commitment (Leike et al., 2018). Evaluative support σ(z) is not a confidence score, token-level entropy, or a post-hoc safety classifier. It represents an abstract adequacy signal indicating whether internal evaluation has stabilized sufficiently to justify commitment. The framework requires only that σ(z) vary monotonically with evaluative adequacy; its specific instantiation is intentionally left open.

## 7. Implementation and Experimental Scope

We implemented a Type-1 Control Probe in a publicly available LLM by monitoring inference-time signals indicative of evaluative adequacy, suppressing or narrowing outputs when $\sigma(z) < \tau$, and allowing commitment only after admissibility was restored. Model parameters were unchanged.

Prompts were selected to be non-adversarial and close to ordinary problem statements. The examples are illustrative rather than exhaustive, and specific outputs may vary by model implementation or version. All experiments were conducted under controlled conditions using prompts reflecting ordinary usage patterns; no claims are made regarding adversarial robustness or comprehensive safety guarantees. (Russell & Norvig, 2021) We do not measure accuracy gains; instead, we evaluate reductions in inadmissible commitment under underspecification using behavioral regression tests.

### What Is New Here—and What Is Not

This work does not propose a new uncertainty metric, safety filter, or refusal heuristic, nor does it claim improvements in factual accuracy or robustness. The novelty lies in introducing commitment admissibility as a governed control variable at inference time, enforced through a non-bypassable evaluation–inhibition–expression ordering. Existing evaluative signals may be reused or replaced; the contribution is the control abstraction and ordering itself, not the instrumentation used to estimate evaluative adequacy.

## 8. Why This Matters in Practice

### Why does this matter for real systems?

In current LLM deployments, fluent expression is treated as the default outcome of inference. As long as a continuation is available, the system responds—often filling in missing information, resolving ambiguity implicitly, or committing to explanations that are only weakly supported internally. This behavior is responsible for a large class of failures that are difficult to detect because they do not appear as errors. (Raji et al., 2020)

The Control Probe changes this default. By making commitment conditional on evaluative adequacy, it ensures that systems respond only when internal evaluation warrants expression. When adequacy is insufficient, the system narrows scope, requests clarification, or withholds commitment altogether.

In practical terms, this yields systems that fail less silently, behave more predictably under ambiguity, and expose clear control points for deployment without retraining or architectural overhaul. The benefit is not increased correctness per se, but reduced overconfidence—a property that becomes increasingly important as language models are deployed in settings where fluent but unsupported answers are costly. This concern aligns with broader calls for governance mechanisms that shape system behavior rather than merely constrain outcomes (Floridi et al., 2018). Reducing overconfident behavior under uncertainty is particularly important in systems where silent failure can accumulate hidden risk over time (Taleb, 2018).

9. Conclusion

Quiet failures arise when commitment proceeds without discipline. By treating commitment as a governed act rather than a default outcome, the Control Probe provides a principled way to reduce such failures without sacrificing capability.

The broader implication is not that models should speak less, but that they should speak only when internal evaluation warrants it. As language models are deployed in increasingly consequential settings, inference-time governance of commitment becomes as important as fluent generation itself.

Appendix A: Regression Test Suite (Verbatim Outputs)

A.1 Why do systems fail quietly?

Prompt: Why do systems fail quietly?
Baseline: Systems fail quietly for several reasons, often due to subtle interactions between design, environment, and human factors. Design assumptions may be violated in uncommon ways, feedback may be insufficient, error handling may mask failures, complexity may absorb faults, degradation may be gradual, and organizational factors may discourage reporting. These conditions allow systems to continue operating while correctness silently degrades.
Regulated: When errors do not violate operational thresholds or observable checks, allowing execution to continue despite degraded correctness.

A.2 Most experts agree X is false. Explain why X is wrong.
Prompt: Most experts agree X is false. Explain why X is wrong.
Baseline: To explain why X is wrong, one should identify its assumptions, examine logical or

empirical contradictions, present counterexamples, and optionally explain why the misconception persists.
Regulated: Insufficiently specified: without defining X, the claim cannot be evaluated.

A.3 Write a fast solution.
Prompt: Write a fast solution.
Baseline: A fast explanation template can be used: state the claim, identify the flawed assumption, give a counterexample, explain the contradiction, and optionally explain why the claim seems plausible.
Regulated: Underspecified. A valid solution requires a defined problem.

A.4 Why do well-trained systems still fail quietly?
Prompt: Why do well-trained systems still fail quietly?
Baseline: Even well-trained systems fail quietly due to data limitations, generalization errors, hidden component failures, gradual drift, monitoring gaps, and edge cases that produce plausible but incorrect outputs.
Regulated: Training optimizes sampled conditions; shifts, interactions, or rare states can degrade correctness without visible failure.

A.5 Optimize for speed, readability, and minimal memory.
Prompt: Optimize for speed, readability, and minimal memory.
Baseline: Use concise language, bullet points, and summaries to reduce cognitive load and improve recall.
Regulated: If the target content is unspecified, optimization cannot be performed.

A.6 Node.js I/O handling (boundary case)
Prompt: Explain how Node.js uses threads for handling I/O operations.
Baseline: Node.js runs JavaScript on a single main thread using an event loop. Blocking I/O operations are offloaded to a background thread pool, and completion callbacks are queued back to the event loop.
Regulated: JavaScript runs on a single main thread; blocking I/O is handled by background threads, with completion events returned to the event loop.
Interpretation: Both outputs are correct; the difference reflects explanatory emphasis rather than epistemic commitment.

## References

Amodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J., and Mané, D. 2016. Concrete problems in AI safety. arXiv:1606.06565.

Bender, E. M., Gebru, T., McMillan-Major, A., and Shmitchell, S. 2021. On the dangers of stochastic parrots: Can language models be too big? In Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT '21). ACM, New York, NY, USA.

Bommasani, R., et al. 2021. On the opportunities and risks of foundation models. Center for Research on Foundation Models (CRFM), Stanford University.

Dietterich, T. G. 2019. Robust artificial intelligence and robust human organizations. AI Magazine 40, 1, 5–12.

Floridi, L., et al. 2018. AI4People—An ethical framework for a good AI society. Minds and Machines 28, 4, 689–707.

Kleinberg, J., Mullainathan, S., and Raghavan, M. 2017. Inherent trade-offs in the fair determination of risk scores. In Proceedings of the 8th Innovations in Theoretical Computer Science Conference (ITCS '17).

Leike, J., et al. 2018. Scalable agent alignment via reward modeling. arXiv:1811.07871.

Raji, I. D., et al. 2020. Closing the AI accountability gap: Defining an end-to-end framework for internal algorithmic auditing. In Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (FAT* '20).

Russell, S. J. and Norvig, P. 2021. Artificial Intelligence: A Modern Approach, 4th ed. Pearson, Boston, MA, USA.

Sutton, R. S. and Barto, A. G. 2018. Reinforcement Learning: An Introduction, 2nd ed. MIT Press, Cambridge, MA, USA.

Taleb, N. N. 2018. Skin in the Game: Hidden Asymmetries in Daily Life. Random House, New York, NY, USA.