

# Soft Constraints Fail Under Pressure: The Case for Hard Admissibility

Constrained decoding, guardrails, and post-generation moderation share one architecture: they shape the model's probability distribution. Under adversarial input, novel context, or distribution shift, soft constraints silently weaken. Hard admissibility, pre-commit, deterministic, typed, does not.

---

## One Architecture Under Many Names

The methods the field reaches for to keep a language model in bounds look diverse and are not. Constrained decoding, logit-biasing re-rankers, classifier-guided generation, constrained beam search, prompt-level guardrails, and post-generation moderation all share a single architecture: they shape the model's probability distribution, or filter its samples, so that disfavored outputs become less likely or are caught after the fact. Some apply a penalty to the logits, some re-rank candidates against a secondary classifier, some wrap the model in an external monitor that inspects inputs and outputs, but the common move is to make the bad output improbable rather than impossible. The constraint rides on top of a statistical process and inherits its statistical nature.

That shared architecture has a shared and predictable weakness. A constraint expressed as a bias on a distribution holds only as well as the distribution behaves, and the distribution stops behaving precisely when it matters most.

## **Why Shaping a Distribution Cannot Hold**

Soft constraints degrade silently under the three conditions that define real deployment. Under adversarial input, a jailbreak, a prompt injection, or a role-play framing reshapes the very distribution the constraint was tuned against, so the penalty that suppressed an output under normal conditions is outweighed under attack, and the system does not announce that its constraint has weakened. Under distribution shift, novel context or an unfamiliar domain moves the model away from the region where the constraint was calibrated, and a bias that was sufficient on the evaluation set is insufficient in the wild. Under scale and composition, constraints interact: a penalty that holds in isolation is overwhelmed when several pressures combine, and the failure surfaces as a quiet leak rather than a visible refusal. In each case the failure mode is the same and is the dangerous one: the constraint does not break loudly, it bends quietly, and the system emits a non-compliant output while reporting nothing wrong.

This is not a tuning problem to be solved with a better classifier or a larger penalty. It is a property of putting the constraint and the thing it constrains in the same statistical layer. A guardrail made of probability can always be outvoted by probability.

## **Hard Admissibility**

The admissibility model places the constraint in a different layer from the generation. The inference engine proposes; the execution substrate disposes. Each candidate transition, including each generation step, is evaluated against typed fields, policy, lineage continuity, entropy bounds, and temporal validity, and the evaluation produces a deterministic outcome of admit, reject, or decompose before the transition is committed. The check is not a bias on the engine's distribution; it is a separate gate the engine's output must pass, and its verdict does not depend on the engine's internal probabilities. An adversarial input can reshape what the engine proposes, but it cannot reshape the gate, because the gate is not a model being prompted, it is a deterministic evaluation over structured state. The failure mode inverts: instead of a silent leak, a

non-compliant proposal produces an explicit rejection recorded in lineage. A constraint that is enforced as a category boundary cannot be made improbable-but-present; it is either satisfied or the step does not happen.

## **Auditability and Regulation**

The difference is decisive for any setting that must be audited. A soft constraint produces no durable account of itself: there is no record of what was suppressed, by how much, or whether the suppression held on a given output, because the constraint lived inside a continuous distribution. Hard admissibility produces exactly that account. Every transition records its admissibility determination in lineage, including the rejections and their reasons, so conformity is demonstrable from the record rather than asserted from a policy document. Regulatory frameworks for high-risk and autonomous systems increasingly demand constraints that the system cannot be argued to have circumvented and evidence that they held; a probabilistic guardrail cannot supply that evidence, and a deterministic, lineage-recorded gate supplies it as a byproduct of operating. This piece is the technical-layer companion to the broader argument that governable behavior comes from architecture rather than from supervision: the reason structure beats supervision here is that supervision is itself a soft constraint, and soft constraints have the failure mode described above.

## **Disclosure Scope**

The structural separation of the inference engine as proposal generator from the execution substrate as authority, and the deterministic admit, reject, or decompose evaluation that treats non-compliant output as a category failure rather than a statistical risk, are disclosed in the cognition filing (U.S. Application No. 19/647,395 and its international counterpart) at Sections 10.5 and 10.8. This article frames those disclosed mechanisms against the constrained-decoding and guardrail families, identifies their common architecture of shaping a probability distribution, and explains

why that architecture degrades silently under adversarial input, distribution shift, and composition while a deterministic admissibility gate does not. References to specific constrained-decoding and guardrail methods are to their public descriptions and are used for comparison only.

---

## **Inference Control** (</inference-control>)

[All 36 steps → \(/inventive-steps\)](/inventive-steps)

Govern inference at the point of generation.

[Explore all disclosures in Inference Control → \(/inference-control\)](/inference-control)

---

[Inference Control overview → \(/inference-control\)](/inference-control)