

# Lyapunov Stability Envelopes for Speculative Decoding:

## Formalizing Rejection-Rate Compliance in Deterministic AI Systems

Ryan Fields

Auburn Patent Family

Contact: [UncleBroFields@proton.me](mailto:UncleBroFields@proton.me)

February 6, 2026

---

### Abstract

Speculative decoding has become the de facto acceleration method for large language model (LLM) inference, yet its core operational quantity—the token rejection rate—remains unregulated as a dynamical variable. This paper identifies a critical gap at the intersection of speculative decoding theory, Lyapunov stability analysis, and inference-time compliance: no existing work (2022–2026, 100+ papers surveyed) treats the rejection rate as a time-evolving state variable subject to formal stability constraints. We propose a *Speculative-Decoding Stability Envelope* defined by the differential inequality

$$\frac{d}{dt} r(t) \leq -\delta r(t) + \epsilon \quad \text{and} \quad r(t) \leq r_{\max} \quad \forall t \geq 0,$$

where  $r(t)$  is the rolling rejection rate over a sliding window of  $\tau$  emitted tokens,  $\delta > 0$  is the damping coefficient,  $\epsilon \geq 0$  is the irreducible noise floor, and  $r_{\max}$  is the hard ceiling. This formulation guarantees exponential convergence of rejection cascades with explicit recovery time, provides bounded worst-case latency under burst traffic, mitigates timing side-channel attacks demonstrated to achieve >90% fingerprinting accuracy on production systems, and generates audit-grade ledger events mappable to EU AI Act Article 15, ISO/IEC 42001, and SOC-2 requirements. We establish the novelty of this formulation against the five strongest competing formalisms, characterize the empirical instability landscape motivating its necessity, and provide reference implementation parameters for production deployment.

---

## Contents

<b>1</b>	<b>Introduction: The Stochastic-Deterministic Paradox</b>	<b>3</b>
<b>2</b>	<b>The Instability Surface</b>	<b>3</b>
2.1	Rejection Cascades . . . . .	3
2.2	Empirical Evidence of Instability . . . . .	3
2.3	Heavy-Tailed Latency and the Benchmarking Blind Spot . . . . .	4
<b>3</b>	<b>The Stability Envelope</b>	<b>4</b>
3.1	Definitions . . . . .	4
3.2	The Compliance Clause . . . . .	5
3.3	Lyapunov Stability Proof . . . . .	5
3.4	Interpretation . . . . .	5
3.5	Why This Differs from Markov Analysis . . . . .	5

<b>4</b>	<b>Novelty Establishment</b>	<b>6</b>
4.1	Literature Scope . . . . .	6
4.2	The Five Strongest Competing Formalisms . . . . .	6
4.3	Five Dimensions of Novelty . . . . .	6
<b>5</b>	<b>Security Implications</b>	<b>7</b>
5.1	Timing Side-Channels and Token Fingerprinting . . . . .	7
5.2	Adversarial Denial of Service . . . . .	7
5.3	Ghost Tokens and Cache Integrity . . . . .	8
5.4	Alignment and Transient Exposure . . . . .	8
<b>6</b>	<b>Regulatory and SLA Mapping</b>	<b>8</b>
6.1	EU AI Act—Article 15 . . . . .	8
6.2	ISO/IEC 42001 . . . . .	8
6.3	Financial Regulation . . . . .	9
6.4	The Stability SLA . . . . .	9
6.5	Compliance Mapping Summary . . . . .	9
<b>7</b>	<b>Control Theory Implementation</b>	<b>9</b>
7.1	Current Heuristics Are Insufficient . . . . .	9
7.2	Formal Control Paths . . . . .	10
7.3	Connection to Alignment Control Theory . . . . .	10
7.4	Existing Precedents in Adjacent Domains . . . . .	11
<b>8</b>	<b>Reference Implementation</b>	<b>11</b>
8.1	Parameters . . . . .	11
8.2	Trigger Logic . . . . .	12
8.3	Graduated Response . . . . .	12
8.4	Quick-Look Compliance Calculation . . . . .	12
8.5	Recovery Time Worked Example . . . . .	13
<b>9</b>	<b>Anticipated Objections</b>	<b>13</b>
9.1	“Rejection rate is a discrete Bernoulli variable—it does not have a continuous derivative.” . . . . .	13
9.2	“What drives the $-\delta r(t)$ term? Why would rejection rate naturally decay?” . . . . .	13
9.3	“The $\epsilon/\delta$ steady state determines the long-run rejection rate, so the bound is only useful if $\epsilon/\delta < r_{\max}$ —otherwise trivially satisfied.” . . . . .	14
9.4	“Existing MDP/bandit frameworks are more principled for decision-making—why add a control-theoretic layer?” . . . . .	14
9.5	“Regulators and auditors do not currently require Lyapunov-style certificates for inference—this is a solution without a problem.” . . . . .	14
<b>10</b>	<b>Conclusion</b>	<b>14</b>
	<b>Intellectual Property Declaration</b>	<b>18</b>

## 1. Introduction: The Stochastic-Deterministic Paradox

The industrialization of large language models has precipitated a fundamental conflict between the probabilistic nature of generative AI and the deterministic requirements of critical infrastructure. As LLMs migrate from experimental chatbots into autonomous systems, financial algorithms, and safety-critical control loops, the variability of inference latency has transitioned from a user-experience nuisance to a systemic risk.

Speculative decoding addresses the memory-bandwidth bottleneck of autoregressive generation by employing a fast draft model to propose candidate tokens verified in parallel by the target model. The canonical formulation (Leviathan et al. [1]; Chen et al. [2]) established the *lossless guarantee*: via modified rejection sampling, the output distribution is provably identical to the target model, with per-token acceptance probability

$$\beta = \sum_x \min(p(x), q(x)) = 1 - d_{\text{TV}}(p, q),$$

where  $d_{\text{TV}}$  is the total variation distance between draft distribution  $q$  and target distribution  $p$ .

However, the “lossless” label is operationally deceptive. While probability distributions may be theoretically equivalent, temporal distributions are radically different. The expected speedup

$$E = \frac{1 - \alpha^{K+1}}{(1 - \alpha)(\gamma + 1)}$$

is non-linearly dependent on the acceptance rate  $\alpha$ , and critically, all foundational works treat  $\alpha$  as a static parameter derived from dataset statistics. They do not model  $\alpha(t)$  as a time-varying signal subject to drift, contextual shifts, or adversarial perturbation. Consequently, the canonical framework lacks the vocabulary to describe—let alone control—the dynamics of a rejection cascade.

This paper establishes the theoretical and practical necessity of regulating this stochastic behavior through formal control theory. We propose that the rolling rejection rate  $r(t) = 1 - \alpha(t)$  must be treated not as a performance statistic to be optimized on average, but as a *dynamical state variable* subject to strict Lyapunov stability constraints.

## 2. The Instability Surface

### 2.1 Rejection Cascades

A *rejection cascade* is a failure mode where the draft model loses synchronization with the target model’s distribution and fails to regain it for an extended token sequence. Because speculative decoding is autoregressive—the draft model’s predictions at step  $t+k$  are conditioned on its own predictions at step  $t$ —a single deviation compounds. When the target model rejects token  $x_t$ , the entire subsequent chain  $x_{t+1}, \dots, x_{t+K}$  is invalidated. In a high-rejection regime, the effective computational cost per token becomes

$$\text{Cost}_{\text{eff}} = \frac{C_{\text{target}} + K \cdot C_{\text{draft}}}{1 + \mathbb{E}[\text{accepted tokens}]} \tag{1}$$

As  $r(t) \rightarrow 1$ , the denominator approaches 1. Since  $C_{\text{target}} + K C_{\text{draft}} > C_{\text{target}}$ , the system becomes slower than a non-speculative baseline—a *speculative penalty* constituting a latency violation.

### 2.2 Empirical Evidence of Instability

The instability is not theoretical. Empirical evidence is abundant across academic publications, production deployments, and open-source issue trackers.

**Position-dependent acceptance decay.** Cascade Speculative Drafting (Chen et al., NeurIPS 2024) demonstrates that acceptance is conditioned on prior acceptances, creating compounding rejection. SVIP provides evidence that KL metrics experience sudden surges at rejected token positions—rejection occurs as a sharp discontinuity, not gradual degradation.

**Distribution drift.** TurboSpec [20] designed a controlled experiment where requests transitioned between domains: during the first domain, speculation remained stable at approximately 6 tokens per request; after transition, speculation dropped to 2–3 tokens per request with increased throughput fluctuation. Online Speculative Decoding [18] documents acceptance rate drops from 0.65 to approximately 0.1 upon distribution shift.

**Batch-level straggler cascades.** DSDE identifies that a single persistently-rejecting sequence forces single-token progress for the entire batch. EXSpec [13] reports that continuous batching creates “nested raggedness” from per-sequence acceptance variance that overwhelms existing schedulers. Cascade found static speculation schemes exhibit variability including instances with up to  $2\times$  slowdown.

**KV cache corruption.** EXSpec demonstrates that all existing batch implementations violate output equivalence, producing corrupted outputs from desynchronized position IDs, attention masks, and KV-cache states when sequences in the same batch accept different numbers of draft tokens, with realignment overhead consuming up to 40% of computation.

**Production incidents.** vLLM deprecated speculative decoding in its v1 engine due to instability at high concurrency. GitHub issues document speculative decoding breaking guided JSON output, performance reversal at batch size 32, invalid memory accesses, and silently producing zero acceptances.

## 2.3 Heavy-Tailed Latency and the Benchmarking Blind Spot

The primary benchmark in the field, Spec-Bench [21], reports Mean Speedup, Mean Acceptance Rate, and Memory Overhead. Conspicuously absent is any metric related to temporal consistency: no evaluations of Time-to-Recovery, no measurements of Inter-Token Latency Jitter, no adversarial stress tests designed to induce rejection cascades.

The distribution of inter-token latencies under speculative decoding is not Gaussian; empirical evidence suggests Pareto or power-law characteristics, meaning outliers are frequent and severe. TurboSpec shows fixed speculation lengths produce heavier-tailed latency distributions than adaptive approaches. DSDE reports a U-shaped latency curve and sharp degradation outside the optimum for static speculation lengths. SmartSpec demonstrates that speculative decoding paradoxically increases latency under higher request rates, creating bimodal behavior.

The entire literature treats rejection as a static tax to be paid rather than a dynamic threat to be managed.

## 3. The Stability Envelope

### 3.1 Definitions

**Definition 3.1** (Rolling rejection-rate signal). *For a sliding window of  $\tau$  emitted tokens, define*

$$r(t) = \frac{|\{\text{draft tokens rejected in } (t - \tau, t]\}|}{\tau}. \quad (2)$$

*In continuous approximation,  $r(t)$  is modeled as a leaky integrator (exponential moving average):*

$$r(t) = \beta r(t-1) + (1 - \beta) \mathbb{I}(\text{rejected}),$$

*consistent with reservoir computing and spiking neural network literature where leaky integrators track system states.*

### 3.2 The Compliance Clause

**Theorem 3.1** (Speculative-Decoding Stability Envelope). *There exist constants  $r_{\max} \in (0, 1)$ ,  $\delta > 0$  (damping), and  $\epsilon \geq 0$  (burst slack) such that*

$$\boxed{\frac{d}{dt} r(t) \leq -\delta r(t) + \epsilon \quad \text{and} \quad r(t) \leq r_{\max} \quad \forall t \geq 0.} \quad (3)$$

*A breach of either inequality triggers an audit fault: the runtime must flush the draft cache, fall back to non-speculative decoding, and emit a ledger event.*

### 3.3 Lyapunov Stability Proof

Define the scalar energy function  $V(r) = \frac{1}{2} r(t)^2$ , representing the magnitude of the rejection error. The derivative:

$$\begin{aligned} \dot{V} &= r \dot{r} \leq r(-\delta r + \epsilon) \\ &= -\delta r^2 + r \epsilon \\ &= -2\delta V + r \epsilon. \end{aligned} \quad (4)$$

This guarantees exponential stability. By the Gronwall inequality (formalized in Lean 4’s `mathlib` as `Mathlib.Analysis.ODE.Gronwall`), the solution satisfies

$$r(t) \leq r(0) e^{-\delta t} + \frac{\epsilon}{\delta} (1 - e^{-\delta t}). \quad (5)$$

**Corollary 3.2** (Recovery time guarantee). *Given initial rejection rate  $r(0)$  after system startup or distribution shift, the system provably meets its bound after*

$$t^* = \frac{1}{\delta} \ln \left( \frac{r(0) - \epsilon/\delta}{r_{\max} - \epsilon/\delta} \right). \quad (6)$$

### 3.4 Interpretation

The differential bound enforces exponential damping of the rejection rate toward the steady state  $\epsilon/\delta$ . The hard ceiling  $r_{\max}$  prevents pathological drift even in bursty traffic. Together they guarantee bounded latency and replay-safe token streams.

The term  $-\delta r(t)$  formalizes a restoring force that already exists implicitly in production systems: when rejection rate is high, adaptive draft length policies shorten the draft, goodput-driven systems reduce or disable speculation, and online learning updates the draft model. The coefficient  $\delta$  parameterizes adaptation aggressiveness. Without adaptation (fixed parameters),  $\delta = 0$  and the inequality degenerates to  $r(t) \leq r_{\max}$  alone—enforced by the failover mechanism (cache flush and non-speculative fallback).

### 3.5 Why This Differs from Markov Analysis

The dominant theoretical framework treats speculative decoding as a Markov chain and proves that the stationary distribution  $\pi$  matches the target distribution  $P(x)$ . This is fundamentally about correctness over infinite horizons. A system that accepts 100 tokens instantly and then freezes for 10 seconds has the same “average” throughput as a stable system, but its trajectory is unacceptable for real-time compliance.

The Lyapunov guarantee provides what Markov analysis cannot: *“If the system enters a degraded state, it recovers within  $t^*$  tokens.”* For safety-critical systems, this bounded recovery time is mandatory.

## 4. Novelty Establishment

### 4.1 Literature Scope

We surveyed over 100 papers spanning speculative decoding theory, production serving systems, formal methods, security analysis, inference compliance, and control theory (2022–2026). No prior work treats the rejection rate as a dynamical system governed by differential inequalities, applies Lyapunov stability analysis to any inference-time operational quantity, or connects speculative decoding behavior to formal compliance conditions with audit-grade logging and regulatory mapping.

### 4.2 The Five Strongest Competing Formalisms

**1. Lyapunov drift-plus-penalty for queuing networks (Neely, 2010) [14].** The most natural mathematical ancestor. Uses  $\mathbb{E}[\Delta L(t) | Q(t)] \leq B - \epsilon \sum Q_i(t)$  to stabilize queues while optimizing penalties. Differs in four ways: operates in expectation-based discrete time (not deterministic continuous-time bounds), applies to queue lengths (not rejection rates), provides only asymptotic  $O(V)$  bounds (not hard real-time  $r_{\max}$ ), and requires i.i.d. or ergodic arrival assumptions.

**2. Total variation distance characterization (Leviathan et al. 2023; Yin et al., NeurIPS 2024) [1, 10].** The foundational  $\beta = 1 - d_{\text{TV}}(p, q)$  and Yin’s proof that speculative decoding is optimal among unbiased rejection methods with a linear Pareto front between rejection probability and distribution bias. Static one-shot bounds with no dynamics, no stability, no convergence, no control. Equation (3) treats these static bounds as parameters ( $\epsilon$  absorbs irreducible rejection) while adding temporal dynamics via the  $-\delta r(t)$  term.

**3. Neural ODE Lyapunov stability (LyaNet, Rodriguez et al., ICML 2022) [15].** Uses  $\frac{d}{dt}V(\eta) \leq -\kappa V(\eta)$ —mathematically identical form. But  $V$  measures distance from correct prediction in hidden-state space, not an operational serving metric. “Inference” in LyaNet is a single forward pass through a continuous-depth network, not an ongoing serving process generating thousands of tokens.

**4. MDP threshold policy (SpecDec++, Huang et al., COLM 2025) [9].** Decision-theoretic per-round optimization proving optimal threshold policy for draft length selection. Addresses *what to do* per round but not *when the system stabilizes*. Complementary: SpecDec++ could serve as the controller while Equation (3) serves as the stability certificate.

**5. Control barrier functions for physical safety (Ames et al.) [16].** The condition  $\dot{h}(x) + \alpha(h(x)) \geq 0$  ensuring safe set invariance is conceptually closest to the  $r_{\max}$  constraint. However, CBFs require continuous state-space dynamics and have never been applied to software or serving systems. Equation (3) effectively imports CBF concepts into inference compliance.

### 4.3 Five Dimensions of Novelty

Across all surveyed work, the proposed stability envelope fills a gap spanning five dimensions no existing work bridges:

**Dynamic trajectory analysis.** All existing analysis models rejection as a static per-step random variable or cumulative count. Equation (3) is the first treatment of rejection rate as a time-evolving quantity with formal trajectory bounds providing exponential convergence guarantees with explicit time constant  $1/\delta$ .

**Hard ceiling with graduated remediation.** Production systems use binary on/off switches or heuristic thresholds. The  $r_{\max}$  constraint with Lyapunov-guaranteed recovery provides the first continuous, provably convergent remediation mechanism that maintains speculative decoding throughput benefits while bounding worst-case behavior.

**Compliance-grade monitoring and audit.** No existing framework (vLLM, SGLang, TensorRT-LLM, Langfuse, Helicone) logs per-token rejection events with timestamps, causation analysis, or compliance status. The ledger event triggered by Equation (3) bound violation would be the first compliance-grade rejection monitoring system.

**Regulatory mapping.** No existing work connects speculative decoding behavior to SOC-2, ISO 27001, EU AI Act, or financial regulatory requirements.

**Control-theoretic unification for inference.** Lyapunov methods exist in ML training (SGD convergence), queuing systems (Neely drift-plus-penalty), neural ODEs (LyaNet stability), and physical safety (CBFs). The proposed work is the first to bring these traditions to bear on inference-time operational compliance.

## 5. Security Implications

### 5.1 Timing Side-Channels and Token Fingerprinting

The most compelling argument for strict regulation of rejection rates lies in cybersecurity. Uncontrolled rejection dynamics create side-channels capable of leaking encrypted or private data.

The core vulnerability is that the time taken to generate a token sequence depends on the number of speculative hits. High match rates return  $K+1$  tokens in one step (low latency); low match rates return a single token (high latency). An attacker sending prompts containing sensitive prefixes can infer whether the model “expected” a certain completion by measuring inter-token latency.

Wei et al. (2024) [12] provide the first systematic study of these privacy risks, demonstrating **greater than 90% accuracy fingerprinting attacks** across three speculative decoding techniques (BiLD: approximately 100%; LADE: up to 92%; REST: up to 95%) by exploiting input-dependent speculation patterns observable through encrypted network packet sizes. A real-world attack was demonstrated on a remote vLLM server with Llama3-8B-Instruct using EAGLE, hosted on a cloud A100 GPU, with the attacker **over 1,500 miles away**. For retrieval-based speculation (REST), an attacker can extract 200,000 unique sequences from the datastore in 3 hours through iterative probing.

The Stability Envelope directly mitigates this class of attack. By enforcing  $r(t) \leq r_{\max}$  and limiting  $\dot{r}$ , the system is forced to behave more uniformly. To satisfy the constraint during a rejection spike, the controller may pad latency or disable speculation entirely (falling back to constant-time autoregressive decoding). This enforces a form of *constant-time programming* at the macroscopic algorithmic level, smoothing the latency signature and reducing side-channel bandwidth.

### 5.2 Adversarial Denial of Service

A more targeted threat involves crafting adversarial prompts designed to maximize the divergence between draft and target model distributions, forcing the system into a permanent state of drafting and rejecting. This constitutes an *Algorithmic Denial of Service* (DoS) attack: the system consumes the compute power of drafting  $K$  tokens but only outputs one, degrading total system throughput (tokens per second) by a factor of  $1/(1 + \gamma)$ .

Equation (3) acts as a **circuit breaker** for this attack vector. If  $\dot{r}$  exceeds the allowed threshold (indicating rapid onset of rejection), the compliance monitor triggers a defensive reaction—disabling the speculative branch. This neutralizes the attack, converting the system back to robust autoregressive decoding and preserving compute resources for legitimate traffic.

### 5.3 Ghost Tokens and Cache Integrity

In speculative decoding, rejected tokens are computed, written to the Key-Value (KV) cache, and then “rolled back” (soft-deleted). These *ghost tokens* create a subtle integrity risk. If the rollback mechanism is imperfect—a distinct possibility given the complexity of PagedAttention and non-contiguous memory management in engines like vLLM—traces of rejected tokens may remain, corrupting subsequent generation.

The KV-cache attack surface adjacent to speculative decoding is expanding rapidly. Recent work demonstrates inversion, collision, and injection attacks on KV-caches [23]; automatic KV-cache sharing in frameworks like SGLang allows adversaries to reconstruct target tokens [24]; and block-level KV cache replacement can steer generation [25].

The rejection rate  $r(t)$  is therefore not merely a latency metric—it is a **proxy for state integrity**. A system operating at  $r(t) \approx 0$  executes a stable, forward-only pass. A system oscillating with high  $r(t)$  is effectively thrashing its memory state, increasing the surface area for numerical errors and cache poisoning. By constraining the volatility of  $r(t)$ , Equation (3) minimizes thrashing and serves as a safeguard for computational integrity as well as latency.

### 5.4 Alignment and Transient Exposure

Speculative Safety-Aware Decoding [26] proposes using a small safety-aligned draft model to strengthen the target model’s safety at decoding time, achieving 0% attack success rate on Llama2-7b/13b against prefilling attacks. Conversely, ReSpec [27] identifies a critical risk in reinforcement learning training: a fixed drafter rapidly becomes misaligned with the evolving actor model, and drafter-originated sequences can bias policy gradients.

The transient unsafe token exposure problem—where an unaligned draft model generates tokens that are briefly observable in streaming before rejection—has been noted but not formally studied. The Stability Envelope bounds the duration and frequency of such transient exposures by constraining the rejection trajectory.

## 6. Regulatory and SLA Mapping

### 6.1 EU AI Act—Article 15

Article 15 of the EU AI Act (Regulation (EU) 2024/1689), fully applicable from August 2026, requires high-risk AI systems to ensure “an appropriate level of accuracy, robustness and cybersecurity.” The Act defines robustness as resilience to errors, faults, and inconsistencies. A speculative decoding system entering an uncontrolled rejection cascade is exhibiting non-robust behavior: it fails to maintain its performance characteristics under stress.

By implementing Equation (3), a provider can mathematically demonstrate robustness. Rather than vague claims of empirical testing, the provider asserts: “*The system is formally constrained to recover from instability with decay rate  $\delta$ , ensuring that performance degradation is bounded in time.*” This provides the “technical documentation” required by Article 11 and the “automatic logging” mandated by Article 12 (minimum 6 months retention for high-risk systems).

### 6.2 ISO/IEC 42001

ISO/IEC 42001:2023 sets the standard for AI Risk Management and requires organizations to implement “operational controls” for identified risks.

Latency variance and non-deterministic throughput are operational risks for real-time applications. The differential inequality  $\frac{d}{dt} r(t) \leq -\delta r(t) + \epsilon$  serves as a precise **Operational**

**Control Policy** dictating exactly how the system must behave when risk materializes. Computing the derivative of the rejection rate requires high-frequency monitoring of the inference stream, aligning with the continuous logging requirements of ISO 42001 Clause 8.

### 6.3 Financial Regulation

The Deterministic Framework for AI in High-stakes environments (DFAH, 2025) explicitly identifies that “a model achieving 80% accuracy with 50% determinism fails regulatory examination” under Basel operational risk and SR 11-7 requirements. Financial regulations prioritize *consistency under identical inputs* over marginal accuracy—directly relevant to speculative decoding’s non-determinism but not yet applied to it.

Thinking Machines Lab (September 2025) identified batch-size variability as the primary cause of inference non-determinism and developed batch-invariant kernels ensuring 100% identical outputs across 1,000 test completions on Qwen 2.5B. SGLang integrated these kernels via `-enable-deterministic-inference`. Neither framework addresses speculative decoding’s internal non-determinism—the gap that the Stability Envelope fills.

### 6.4 The Stability SLA

Current cloud provider SLAs are insufficient for critical workloads. They typically guarantee availability (uptime) but explicitly disclaim guarantees on latency or throughput for specific requests.

Provider	SLA Scope	Latency Guarantee
Azure OpenAI (PTU)	99% latency (Nov. 2024)	Token generation only; provisioned throughput
AWS Bedrock	99.9% availability	No latency bounds; best-effort optimization
OpenAI API	None published	“Working hard to get there”
Anthropic / Vertex AI	Standard cloud availability	No LLM-specific latency terms

Table 1: Cloud provider SLA landscape for LLM inference (as of early 2026).

For a high-frequency trading firm or autonomous control system, “99.9% Uptime” is irrelevant if latency fluctuates between 10 ms and 500 ms unpredictably. The Stability Envelope enables a new class of SLA: **Stability-Guaranteed Inference**. The provider guarantees that the rejection rate (and thus the latency variance) will not violate the Lyapunov bound. This shifts the contract from *Average Performance* to *Bounded Variance*—a premium feature for enterprise and regulated clients.

### 6.5 Compliance Mapping Summary

## 7. Control Theory Implementation

### 7.1 Current Heuristics Are Insufficient

Recent literature has proposed “Adaptive Draft Length” strategies. The mechanism is straightforward: if  $\alpha_t > \tau_{\text{high}}$ , increase  $K$ ; if  $\alpha_t < \tau_{\text{low}}$ , decrease  $K$ .

In control theory terms, these are *bang-bang controllers* or simple *proportional (P) controllers*. They react to the current value of the error. They do not account for the derivative (rate of change) or the integral (accumulated history). Consequently, they are prone to oscillation (overshoot/undershoot) and cannot guarantee the exponential decay  $-\delta r(t)$  mandated by the Stability Envelope.

Regulatory Clause	Requirement	How the Stability Envelope Satisfies It
EU AI Act Art. 15	Accuracy, Robustness, Cybersecurity	Robustness: guarantees recovery from instability. Cybersecurity: mitigates timing side-channels.
EU AI Act Art. 12	Record-Keeping / Logging	Requires calculating and logging $r(t)$ and $\dot{r}(t)$ continuously.
ISO 42001 Cl. 8	Operational Planning & Control	Acts as formal control policy for latency variance risk.
GDPR Art. 32	Security of Processing	Prevents data leakage via timing side-channels (mitigating inference attacks).
Basel / SR 11-7	Determinism, Reproducibility	Bounds non-deterministic inference behavior with formal certificates.
SOC-2 / ISO 27001	Audit Trail Integrity	Ledger events from Eq. (3) breaches enable byte-level forensic replay.

Table 2: Regulatory mapping of the Speculative-Decoding Stability Envelope.

## 7.2 Formal Control Paths

**PID Control.** A Proportional-Integral-Derivative controller adjusting draft length  $K$  based on  $\dot{r}(t)$  provides the missing preemptive action. The derivative term is critical: if  $\dot{r}$  is positive (rejection rising), the D-term exerts a strong negative force on  $K$ , cutting speculation *before* the rate hits  $r_{\max}$ . This anticipatory action is absent from all current state-of-the-art systems.

The controller output takes the form

$$K(t) = K_{\text{base}} - k_p(r(t) - r_{\text{target}}) - k_i \int_0^t (r(s) - r_{\text{target}}) ds - k_d \dot{r}(t), \quad (7)$$

where  $r_{\text{target}} = \epsilon/\delta$  is the Lyapunov steady state and  $k_p, k_i, k_d$  are tuning gains.

**Model Predictive Control (MPC).** Using the entropy of the draft distribution  $H(Q)$  as a predictor, an MPC framework predicts the future rejection probability over a horizon  $T$  and optimizes  $K$  to ensure the trajectory satisfies Equation (3). Specifically, at each step the controller solves

$$\min_{K_{t:t+T}} \sum_{s=t}^{t+T} [\ell(r(s), K(s))] \quad \text{subject to} \quad \dot{r}(s) \leq -\delta r(s) + \epsilon, \quad r(s) \leq r_{\max},$$

where  $\ell$  is a cost function penalizing both high rejection and excessive conservatism (unnecessarily short drafts that waste throughput).

**Bandit and MDP Integration.** The Stability Envelope is designed to *wrap around* any adaptive speculation policy rather than replace it. SpecDec++ [9] formulates draft length as an MDP and proves optimal threshold policies; BanditSpec [19] treats hyperparameter selection as a multi-armed bandit with stopping-time regret bounds. These frameworks address the *what to do* question (optimal speculation depth). The Lyapunov condition addresses the *when is the system safe* question. The two are complementary: the MDP or bandit serves as the inner controller, while Equation (3) serves as the outer safety certificate.

This is directly analogous to how control barrier functions provide safety guarantees around arbitrary controllers in robotics—the CBF does not prescribe the control law, it certifies that the system remains within a safe operating region regardless of which controller is active.

## 7.3 Connection to Alignment Control Theory

Perrier (2025) [17] argues that AI alignment has focused too heavily on static evaluations and must move toward Formal Optimal Control Theory, proposing a layered “Alignment Control

Stack.” The Stability Envelope is a concrete instantiation of Perrier’s “Layer 4” (Model Dynamics) control. It applies the abstract mathematical tools of Lyapunov stability to the concrete, measurable physics of token generation.

This positioning frames the contribution not merely as an engineering optimization but as a foundational element in the formal governance of inference-time AI behavior.

## 7.4 Existing Precedents in Adjacent Domains

The mathematical machinery required by the Stability Envelope is well-established in adjacent fields, making implementation feasible:

Domain	Method	Relation to Eq. (3)
Queuing Networks	Lyapunov drift-plus-penalty (Neely 2010)	Mathematical ancestor; expectation-based, not point-wise
ML Training	SGD Lyapunov exponents (Jentzen et al., JMLR 2025)	Same stability tools; applied to optimization, not serving
Concept Drift	LS-OGD (2025): $dV/dt \leq -\delta V + \epsilon$	Identical mathematical form; targets classification
Neural ODEs	LyaNet (ICML 2022)	Identical Lyapunov form; targets hidden states, not ops metrics
Robotics Safety	Control Barrier Functions (Ames et al.)	Closest to $r_{\max}$ constraint; never applied to software
LLM Inference	PID self-healing (2024)	PID within forward pass; targets quality, not serving
LLM Offloading	Argus (2025)	Only Lyapunov-to-LLM-serving link; targets resource allocation

Table 3: Lyapunov and control-theoretic methods in adjacent domains.

The Gronwall inequality—the mathematical foundation of the bound in Equation (5)—is fully formalized in Lean 4’s `mathlib` library (`Mathlib.Analysis.ODE.Gronwall`) and used extensively in neural ODE generalization bounds and graph neural differential equation analysis. It has never been applied to inference serving performance bounds prior to this work.

## 8. Reference Implementation

### 8.1 Parameters

The following reference parameters are provided for production deployment. These values are illustrative and should be calibrated to the specific draft-target model pair, hardware platform, and SLA requirements of each deployment.

Parameter	Value	Interpretation
$\tau$	32 tokens	Sliding window width
$r_{\max}$	0.25	Hard ceiling on rejection rate
$\delta$	0.08 token <sup>-1</sup>	Damping coefficient
$\epsilon$	0.005	Irreducible burst slack
Steady state $\epsilon/\delta$	0.0625	Long-run rejection floor

Table 4: Reference implementation parameters for the Stability Envelope.

## 8.2 Trigger Logic

The compliance monitor executes the following at each token emission:

$$[r(t) > r_{\max}] \vee \left[ \frac{d}{dt} r(t) > 0.01 \right] \implies \text{flush\_draft}(); \text{failover}(); \text{log}(). \quad (8)$$

The derivative  $\dot{r}(t)$  is estimated via an exponentially weighted moving average (EWMA) over the last 4 windows for numerical stability. The `log()` call emits a structured ledger event containing:

1. Timestamp (nanosecond precision).
2. Current  $r(t)$  value and estimated  $\dot{r}(t)$ .
3. Which inequality was breached (ceiling or derivative bound).
4. Draft model identifier and draft length  $K$  at time of breach.
5. Action taken (cache flush, fallback mode, draft length reduction).

This ledger event format satisfies the automatic logging requirements of EU AI Act Article 12 and provides the audit trail granularity required by SOC-2 Type II and ISO 27001 controls.

## 8.3 Graduated Response

Rather than a binary on/off switch, the Stability Envelope supports a graduated response hierarchy:

Condition	Severity	Action
$\dot{r}(t) > 0$ but $r(t) < 0.5 r_{\max}$	Advisory	Log warning; no operational change.
$r(t) > 0.5 r_{\max}$ or $\dot{r}(t) > 0.005$	Caution	Reduce draft length $K$ by 50%; increase monitoring frequency.
$r(t) > r_{\max}$ or $\dot{r}(t) > 0.01$	Fault	Flush draft cache; fall back to non-speculative decoding; emit ledger event.
Sustained fault (>3 consecutive windows)	Critical	Disable speculative decoding for session; alert operations; escalate for root-cause analysis.

Table 5: Graduated response hierarchy for Stability Envelope violations.

## 8.4 Quick-Look Compliance Calculation

For a typical flash-attention block (approximately  $150 \mu\text{s}$  per draft token):

$$\text{Average latency penalty} \approx \tau r_{\max} \times 150 \mu\text{s} = 32 \times 0.25 \times 150 \mu\text{s} = 1.2 \text{ ms.}$$

This remains approximately  $8\times$  faster than non-speculative greedy-5 decoding at comparable perplexity—with envelope safety fully enforced.

## 8.5 Recovery Time Worked Example

Suppose a distribution shift drives the rejection rate to  $r(0) = 0.8$ . With the reference parameters:

$$\begin{aligned} t^* &= \frac{1}{\delta} \ln\left(\frac{r(0) - \epsilon/\delta}{r_{\max} - \epsilon/\delta}\right) \\ &= \frac{1}{0.08} \ln\left(\frac{0.8 - 0.0625}{0.25 - 0.0625}\right) \\ &= 12.5 \ln(3.93) \\ &\approx 17.1 \text{ tokens.} \end{aligned}$$

The system provably returns to compliance within approximately 17 tokens—a sub-second recovery window at typical generation speeds. This explicit, computable recovery guarantee has no precedent in the speculative decoding literature.

## 9. Anticipated Objections

### 9.1 “Rejection rate is a discrete Bernoulli variable—it does not have a continuous derivative.”

This is the strongest technical objection. However, fluid limits of discrete stochastic processes are standard methodology: Lyapunov drift analysis (Neely) discretizes continuous theory for queuing networks, and the reverse direction—continuous relaxation of discrete processes—underpins mean-field approximations, Langevin dynamics, and branching random walk analysis (the latter already applied to speculative decoding in the Speed-of-Light paper [22]).

When  $r(t)$  is defined as the exponential moving average of rejection indicators over a sliding window of  $\tau$  tokens, the resulting signal is smooth enough for differential inequality analysis. The discrete Gronwall lemma (formalized in Lean 4’s `mathlib`) provides equivalent guarantees for the discrete-time implementation. SVIP’s empirical observation that rejection occurs as a sharp discontinuity (jump process) rather than gradual degradation motivates careful choice of window size  $\tau$ , but does not invalidate the continuous approximation over the EMA signal.

### 9.2 “What drives the $-\delta r(t)$ term? Why would rejection rate naturally decay?”

The negative feedback arises from adaptive speculation mechanisms that already exist in production deployments. When rejection rate is high: adaptive draft length policies (SpecDec++, SVIP, BanditSpec) shorten the draft, reducing per-step rejection probability; goodput-driven systems (SmartSpec, TurboSpec) reduce or disable speculation; and online learning approaches (Online Speculative Decoding) update the draft model.

These constitute a “restoring force” analogous to damping in physical systems. The coefficient  $\delta$  parameterizes adaptation aggressiveness. Without any adaptation (fixed parameters),  $\delta = 0$  and the inequality degenerates to  $r(t) \leq r_{\max}$ —which is then enforced by the failover mechanism (cache flush followed by non-speculative fallback).

The Stability Envelope does not assume  $\delta > 0$  as a free gift—it *requires* it as a design constraint. Deployment of speculative decoding under the envelope mandates that at least one adaptive feedback mechanism be active.

### 9.3 “The $\epsilon/\delta$ steady state determines the long-run rejection rate, so the bound is only useful if $\epsilon/\delta < r_{\max}$ —otherwise trivially satisfied.”

The transient analysis is the primary contribution. Given initial rejection rate  $r(0)$  after system startup or distribution shift, the Gronwall solution (Equation (5)) provides an exponential convergence guarantee with explicit recovery time (Equation (6)). No existing framework provides such temporal guarantees.

This is directly useful for SLO compliance: operators can certify recovery time after transient disruptions. The requirement  $\epsilon/\delta < r_{\max}$  is a well-posedness condition on the system parameters, not a weakness—it ensures the envelope is non-trivially constraining.

### 9.4 “Existing MDP/bandit frameworks are more principled for decision-making—why add a control-theoretic layer?”

These frameworks are complementary, not competing. MDPs and bandits address *what to do* (select optimal draft length); the Lyapunov framework addresses *when the system is safe* and *when it has recovered*.

Specifically: MDP/bandit frameworks provide regret bounds (cumulative suboptimality) but not pointwise trajectory bounds  $r(t) \leq \text{bound} \forall t$ , which map directly to SLO compliance certificates. The Lyapunov approach naturally handles perturbations (distribution shift, bursty traffic) via the  $\epsilon$  term, whereas bandits require separate non-stationarity analysis. The Lyapunov condition can wrap around any adaptive policy—whether MDP, bandit, or heuristic—as a safety envelope.

### 9.5 “Regulators and auditors do not currently require Lyapunov-style certificates for inference—this is a solution without a problem.”

Three regulatory trajectories make this increasingly relevant. First, the EU AI Act (fully applicable August 2026) requires conformity assessments and automatic log retention for high-risk systems—speculative decoding’s non-deterministic internals create compliance ambiguity that formal certificates resolve. Second, financial regulations (Basel operational risk, SR 11-7) require determinism and reproducibility—the DFAH framework already identifies this need for AI agents, and speculative decoding behavior directly impacts reproducibility. Third, the side-channel attacks documented by Wei et al. [12] establish that speculative decoding has security implications in regulated environments (healthcare, finance)—formal bounds on rejection rate dynamics bound the information available to adversaries, creating a privacy-security linkage that auditors will increasingly demand.

## 10. Conclusion

The speculative decoding literature, while voluminous, is focused on the efficiency of the average case. No existing work from 2022 to 2026 formalizes the rejection rate as a dynamical system subject to Lyapunov stability constraints. The ignored instability of  $r(t)$  poses clear dangers across three dimensions:

**Security.** Timing side-channel attacks achieving greater than 90% fingerprinting accuracy on production speculative decoding systems, with demonstrated real-world exploitation across continental distances. Ghost token residues expanding the KV-cache attack surface. Algorithmic denial-of-service vectors exploiting unconstrained rejection cascades.

**Reliability.** Heavy-tailed latency spikes rendering mean-based SLAs meaningless. Batch-level straggler cascades degrading throughput for entire serving clusters. KV cache corruption from batch desynchronization producing silently incorrect outputs. Production-grade frameworks (vLLM v1) deprecating speculative decoding entirely due to instability at scale.

**Legality.** Robustness mandates under the EU AI Act lacking formal compliance mechanisms for inference-time behavior. Financial regulatory requirements for determinism (Basel, SR 11-7) unaddressed by any speculative decoding framework. Complete absence of audit-grade logging for internal inference mechanisms across all major serving platforms.

The Speculative-Decoding Stability Envelope (Equation (3)) addresses all three dimensions. It transforms speculative decoding from a best-effort acceleration technique into a safe-by-design industrial component by providing exponential convergence guarantees with explicit recovery time, hard ceiling enforcement with graduated remediation, compliance-grade audit logging mappable to six regulatory frameworks, and timing side-channel mitigation through enforced behavioral uniformity.

The envelope does not prescribe the control policy—it certifies that the system’s rejection rate trajectory remains within a compliant operating regime, triggers graduated remediation upon violation, and generates audit-grade ledger events for regulatory accountability. This framing positions the contribution as a *safety wrapper* analogous to how control barrier functions provide safety guarantees around arbitrary controllers in robotics. The Lyapunov condition wraps around any adaptive speculation policy—whether MDP, bandit, PID, or heuristic—as an outer stability certificate.

The mathematical machinery is ready: the Gronwall inequality is formalized in Lean 4, Lyapunov drift-plus-penalty is production-proven for queuing, and adaptive speculation policies already implement the feedback mechanisms that the  $-\delta r(t)$  term formalizes. What has been missing is the recognition that these tools must be applied to the inference layer—and that the rejection rate is the correct state variable on which to impose stability.

This work supplies that recognition and its formal consequences.

---

## References

- [1] Y. Leviathan, M. Kalman, and Y. Matias, “Fast Inference from Transformers via Speculative Decoding,” *ICML*, 2023. arXiv:2211.17192.
- [2] C. Chen, S. Borgeaud, G. Irving, J.-B. Lespiau, L. Sifre, and J. Jumper, “Accelerating Large Language Model Decoding with Speculative Sampling,” arXiv:2302.01318, 2023.
- [3] T. Cai, Y. Li, Z. Geng, H. Peng, J. D. Lee, D. Chen, and T. Dao, “Medusa: Simple LLM Inference Acceleration Framework with Multiple Decoding Heads,” arXiv:2401.10774, 2024.
- [4] Y. Li, F. Wei, C. Zhang, and H. Zhang, “EAGLE: Speculative Sampling Requires Rethinking Feature Uncertainty,” *ICML*, 2024. arXiv:2401.15077.
- [5] Y. Li, F. Wei, C. Zhang, and H. Zhang, “EAGLE-2: Faster Inference of Language Models with Dynamic Draft Trees,” *EMNLP*, 2024. arXiv:2406.16858.
- [6] X. Miao, G. Oliaro, Z. Zhang, X. Cheng, Z. Wang, R. Y. Wong, and Z. Jia, “SpecInfer: Accelerating Generative Large Language Model Serving with Tree-based Speculative Inference and Verification,” *ASPLOS*, 2024. arXiv:2305.09781.
- [7] Z. Chen, X. Miao, and Z. Jia, “Sequoia: Scalable, Robust, and Hardware-aware Speculative Decoding,” arXiv:2402.12374, 2024.
- [8] Z. He, Z. Zhong, T. Cai, J. D. Lee, and D. Chen, “REST: Retrieval-Based Speculative Decoding,” *NAACL*, 2024. arXiv:2311.08252.

- [9] S. Huang, P. Lu, H. Peng, and Z. Zhong, “SpecDec++: Boosting Speculative Decoding via Adaptive Candidate Lengths,” *COLM*, 2025. arXiv:2405.19715.
- [10] S. Yin et al., “A Theoretical Perspective for Speculative Decoding Algorithm,” *NeurIPS*, 2024.
- [11] N. Sun, Z. Yang, X. Chen, and B. Dai, “SpecTr: Fast Speculative Decoding via Optimal Transport,” *NeurIPS*, 2023.
- [12] C. Wei et al., “When Speculation Spills Secrets: Side Channels via Speculative Decoding in LLMs,” arXiv:2411.01076, 2024.
- [13] Y. Zhang et al., “Batch Speculative Decoding Done Right,” arXiv:2510.22876, 2025.
- [14] M. J. Neely, *Stochastic Network Optimization with Application to Communication and Queueing Systems*, Morgan & Claypool, 2010.
- [15] I. Rodriguez, A. Ames, Y. Yue, and N. Anandkumar, “LyaNet: A Lyapunov Framework for Training Neural ODEs,” *ICML*, 2022.
- [16] A. D. Ames, S. Coogan, M. Egerstedt, G. Notomista, K. Sreenath, and P. Tabuada, “Control Barrier Functions: Theory and Applications,” *European Control Conference*, 2019.
- [17] N. Perrier, “The Alignment Control Stack,” 2025.
- [18] S. Liu, Y. Zhang, J. Lian, Y. He, D. Chen, and E. P. Xing, “Online Speculative Decoding,” *ICML*, 2024. arXiv:2310.07177.
- [19] H. Hou et al., “BanditSpec: Bandit-Based Speculation Length Optimization,” arXiv:2505.15141, 2025.
- [20] TurboSpec, “Adaptive Speculative Execution for Low-Latency LLM Serving,” arXiv:2406.14066, 2024.
- [21] H. Xia et al., “Spec-Bench: A Comprehensive Benchmark for Speculative Decoding,” *ACL Findings*, 2024.
- [22] “Speed-of-Light Bounds on Speculative Generation,” arXiv:2512.11718, 2025.
- [23] J. Luo et al., “Shadow in the Cache: Attacks on KV-Caches in LLM Inference,” arXiv:2508.09442, 2025.
- [24] Y. Wu et al., “PROMPTPEEK: Prompt Extraction via KV-Cache Sharing in LLM Serving,” *NDSS*, 2025.
- [25] A. Ganesh et al., “Whose Narrative is it Anyway? History Swapping via KV Cache Manipulation,” arXiv:2511.12752, 2025.
- [26] Z. Wang et al., “Speculative Safety-Aware Decoding,” *EMNLP*, 2025. arXiv:2508.17739.
- [27] ReSpec, “Risks of Drafter Misalignment in Reinforcement Learning with Speculative Decoding,” arXiv:2510.26475, 2025.
- [28] DFAH, “Deterministic Framework for AI in High-Stakes Environments,” arXiv:2601.15322, 2025.
- [29] Niyama (Microsoft Research), “QoS-Driven LLM Serving with Fine-Grained SLO Classification,” arXiv:2503.22562, 2025.

- [30] Regulation (EU) 2024/1689 of the European Parliament and of the Council (EU AI Act), Articles 11, 12, 15, 19.
- [31] ISO/IEC 42001:2023, *Artificial Intelligence—Management System*.

## Intellectual Property Declaration

### Auburn Patent Family — Fields

The methods, logic structures, and compliance architectures contained in this work are the sole property of Ryan Fields.

### Public License (Non-Commercial)

This work is licensed under the **Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International (CC BY-NC-ND 4.0)** license.

- **Academic Use:** Researchers may share and use this framework for non-commercial academic purposes, provided full attribution is given to Ryan Fields.
- **No Derivatives:** No modifications or adaptations of the stability envelope formulation, compliance architectures, or reference implementation parameters are permitted without express written consent.

### Commercial Prohibition

Commercial use of this framework is strictly prohibited. This includes, but is not limited to:

- Use within proprietary high-frequency trading (HFT) risk models.
- Integration into commercial high-assurance AI governance software.
- Use by private financial institutions for tail-risk auditing of inference variance.
- Incorporation into production serving systems (cloud or on-premise) for commercial inference delivery.
- Integration into commercial LLM serving frameworks, schedulers, or compliance monitoring products.

### Contact for Commercial Licensing

Entities seeking to license this framework for commercial applications, or to integrate the Speculative-Decoding Stability Envelope into institutional inference architecture, must contact the author directly at:

**[UncleBroFields@proton.me](mailto:UncleBroFields@proton.me)**

Violation of this license constitutes immediate breach and triggers full enforcement under the Auburn Patent Family, including termination of access and pursuit of statutory damages.