



[Home](#) [Licensing](#) [Patents](#) [Articles](#)

Together AI Optimizes Inference Speed, Not Inference Governance

by [Nick Clark](#) | Published March 27, 2026 | [PDF](#)

Together AI built a high-performance inference platform that serves open-source models at competitive speed and cost. The infrastructure engineering to achieve fast inference across diverse model architectures is substantial. But Together AI's platform optimizes the delivery of model output without evaluating that output's semantic admissibility. The model generates, the infrastructure serves it fast, and the application receives it. Inference control provides the structural gate that evaluates output against persistent semantic state at the point of generation, without sacrificing the throughput that makes the platform valuable.

What Together AI built

Together AI provides inference APIs for a wide range of open-source models with optimized serving infrastructure. Custom hardware, efficient batching, and model-specific optimizations achieve low latency and high throughput. The platform supports chat, completion, embedding, and image generation endpoints. Fine-tuning services allow customization. The value proposition is clear: fast, affordable inference for open-source models without managing infrastructure.

Output governance at the platform level includes basic safety filtering. Beyond this, the output is the model's output, served at maximum speed. The platform's competitive advantage is throughput and cost, not output governance. This is an appropriate engineering priority for infrastructure, but it leaves governance entirely to the caller.

The gap between fast and governed

Speed and governance are not inherently in tension. An admissibility gate that evaluates output against persistent state can operate within the serving pipeline without adding significant latency if it is architecturally integrated. The gate evaluates the semantic properties of the output against pre-registered constraints. This evaluation is lightweight compared to model inference itself. The latency cost is minimal. The governance value is structural.

Without this gate, every application using Together AI's inference must independently implement output governance. The result is inconsistent governance quality across applications, with some implementing robust evaluation and many implementing none. The applications that most need governance, those built by smaller teams with limited AI safety expertise, are least likely to implement it independently.

What inference control enables

With an admissibility gate integrated into the serving pipeline, Together AI's platform provides governed inference as an infrastructure capability. Applications register semantic constraints alongside their API configuration. Every output is evaluated against these constraints before returning. The gate adds minimal latency while providing structural assurance that output is semantically admissible.

The model-agnostic property means the same governance mechanism works across all models served by the platform. Whether the application uses Llama, Mistral, or any other hosted model, the admissibility gate evaluates output against the same application-specific semantic constraints. This gives applications consistent governance regardless of which model they choose.

The structural requirement

Together AI's inference performance is competitive. The gap is governing what is served, not just serving it fast. Inference control provides the admissibility gate that can be integrated into high-throughput serving pipelines with minimal latency impact while ensuring that every output is evaluated against persistent semantic state. The inference platform that governs output is a more complete infrastructure offering than one that serves output at maximum speed without evaluation.

[Inference Control All 21 steps →](#)

Govern inference at the point of generation.

Primary Technical Disclosure

[◦ Inference-Time Semantic Execution Control](#)

Secondary Technical

[◦ Inference as Semantic Execution](#)◦ [Semantic Admissibility Gate](#)◦ [Entropy-Bounded Semantic Admissibility](#)◦ [Inference-Time Semantic Budget](#)◦ [Semantic Rollback and Checkpoint Recovery](#)◦ [Multi-Model Arbitration With Shared Semantic State](#)◦ [Structural Elegance Evaluation](#)◦ [Rights-Grade Inference Governance](#)◦ [Semantic State Object](#)◦ [Semantic State Object Schema](#)◦ [Inference Transition as Mutation](#)◦ [Trust-Slope Continuity Across Inference](#)◦ [Anchored Semantic Resolution](#)◦ [Semantic Lineage Recording](#)◦ [Policy-Governed Inference Execution](#)◦ [Partial State Handling](#)◦ [Model-Agnostic Inference Governance](#)◦ [Pre-Generation vs Post-Generation Distinction](#)◦ [Affect-Modulated Inference Admissibility](#)◦ [Integrity-Aware Inference](#)◦ [Confidence-Gated Inference Advancement](#)◦ [Inference Deployment Embodiments](#)

Applications (General)

[◦ Safety Without Alignment Theater: Why Structure Beats Supervision](#)◦ [How Commercial AI Platforms Reduce Prompt Size, Drift, and Governance Risk at Scale](#)◦ [When Execution Governance Becomes a Competitive Advantage — The Layer After LLM Gateways](#)◦ [Enterprise LLM Governance at the Point of Generation](#)◦ [Healthcare AI Admissibility Before Clinical Output](#)◦ [Inference Control for Legal Document Generation](#)◦ [Inference Control for Financial Advisory Output](#)◦ [Inference Control for Education Content Generation](#)◦ [Inference Control for Government Communications](#)

Applications (Specific)

[◦ Einstein Generates Without Semantic Admissibility](#)◦ [Databricks Serves Inference Without Semantic Gates](#)◦ [Snowflake Cortex Generates Without Admissibility Gates](#)◦ [Hugging Face Serves Models Without Semantic Governance](#)◦ [Cohere's Enterprise LLM Has No Semantic Admissibility Gate](#)◦ [Together AI Optimizes Inference Speed, Not Inference Governance](#)◦ [SageMaker Serves Models Without Semantic Admissibility](#)◦ [Vertex AI Generates Without Per-Transition Admissibility](#)◦ [Azure ML Deploys Models Without Admissibility Gates](#)◦ [Modal Runs Inference Fast Without Governing Output](#)◦ [Replicate Serves Open Models Without Semantic Governance](#)◦ [Fireworks AI Optimizes Speed Without Governing Semantics](#)◦ [Groq's LPU Accelerates Inference Without Governing It](#)◦ [Cerebras Achieves Wafer-Scale Inference Without Semantic Governance](#)

[Inference Control overview →](#)

AQ

deterministic

autonomy

Legal

Subject to one or more pending U.S. and international patent applications, see [Patents](#) for the current list and status. No license, express or implied, is granted. Any use requires a separate written agreement—see [Licensing](#). Patent applications referenced on this site are pending. Claim scope, if any, is subject to examination and may issue in altered form or not at all. See [Legal](#) for terms and conditions.

Adaptive Query™ is a trademark of Nicholas Clark. U.S. federal registration is pending. federal registration. AQ™, AQ Inside™, Adaptive Index™, Adaptive Network™, Semantic Agent™, @AQ™, AQID™, and Adaptive Coin™ are used as trademarks in connection with the Adaptive Query platform

and brand. Other names may be trademarks of their respective owners.

Platform operated by Adaptive Query LLC, which provides patent and trademark licensing services. Copyright © 2025-2026 Nicholas Clark. All rights reserved.

Last updated: 2026-03-03



- [Inventive Steps](#)
- [Licensing](#)
- [Patents](#)
- [Articles](#)
- [Legal](#)
- [Opportunities](#)
- [Sitemap](#)



-
- nick@qu3ry.net
- 72 28 14 36 01



[Invented by Nick Clark](#) | Founding Investors: Devin Wilkie