

Volume 2, Issue 1

Research Article

Date of Submission: 16 Mar, 2026

Date of Acceptance: 13 Apr, 2026

Date of Publication: 20 Apr, 2026

Constraint Topology and the Economics of Deterministic AI Infrastructure a Four-Run NAGI Feasibility Study of Agentic, DeClawed, Hybrid, and Real-World AI Deployment Architectures

Ean Mikale*

Founder & Chief Executive Officer, Infinite 8 Industries, Inc., USA

*Corresponding Author: Ean Mikale, Founder & Chief Executive Officer, Infinite 8 Industries, Inc., USA.

Citation: Mikale, E. (2026). Constraint Topology and the Economics of Deterministic AI Infrastructure a Four-Run NAGI Feasibility Study of Agentic, DeClawed, Hybrid, and Real-World AI Deployment Architectures. *Adv Brain-Computer Interfaces Neural Integr*, 2(1), 01-09.

Abstract

We present a four-run empirical study of AI infrastructure cost economics using the NAGI (Non-Agentic General Intelligence) Feasibility Engine — a deterministic constraint processing system that identified and certified mathematically valid operating configurations for AI deployments. Across four sequential constraint runs spanning 10 to 17 operational parameters, 2,000 sampled states per run, and four distinct architectural scenarios (Agentic Baseline, DeClawed, Hybrid Enterprise, and Real-World Pricing), we demonstrate that AI infrastructure economics is expressible as a fully solvable constraint optimization problem with stable topology. Our central empirical finding is a constraint regime shift: agentic architectures are bound by financial constraints (infrastructure cost saturation at 89–92%), while deterministic architectures are bound by physical efficiency constraints (PUE and token throughput at 91–93%), and hybrid architectures are governed by traffic allocation identity constraints (94% saturation). We further show that under real-world market pricing conditions (AWS GPU rates \$2–\$20/hr, input token costs \$1.00–\$1.80/million, energy \$60–\$75/MWh), deterministic architectures maintain linear cost scaling while agentic systems exhibit superlinear cost growth [1–6]. The feasibility boundary exhibits a Hausdorff dimension of 1.000 (smooth) and entropy stability of $HR = 3.9120$ across all runs, certifying the absence of hidden instabilities or discontinuities. We argue that the perceived unpredictability of AI systems is not an intrinsic property of intelligence, but a consequence of unconstrained architecture.

Keywords: AI Infrastructure Economics, Constraint Optimization, Deterministic AI, Agentic Systems, NAGI, Feasibility Analysis, Cost Scaling, Latin Hypercube Sampling

Classification: Operational Research

AI Systems Engineering

Infrastructure Economics

Introduction

The economics of AI infrastructure are widely described as unpredictable. Organizations deploying large language models at scale routinely report cost overruns, utilization failures, and nonlinear demand spikes that resist conventional capacity planning methods [1]. The dominant explanation attributes this unpredictability to the inherent stochasticity of large generative models: token generation is probabilistic, agent orchestration is recursive, and demand patterns are non-stationary [7,8].

This paper challenges that explanation. We argue that the perceived unpredictability of AI infrastructure costs is not an intrinsic property of the underlying intelligence, but a consequence of unconstrained architecture. When AI systems are deployed without rigorous constraint specification, their economic behavior appears stochastic because no feasibility envelope has been defined. The system is not unpredictable — it is underdetermined.

We test this argument using the NAGI Feasibility Engine, a deterministic constraint processing system developed by Infinite 8 Industries, Inc. NAGI accepts a set of operational constraints and processes them via Latin Hypercube Sampling (LHS) across a user-specified state space, identifying all mathematically valid operating configurations simultaneously

[9]. Unlike simulation or Monte Carlo methods, NAGI produces exact solutions to the feasibility problem: it certifies which combinations of operational parameters simultaneously satisfy all constraints, and classifies the topology of the resulting feasible region [10-12].

Contribution

This paper makes four primary contributions:

- We formalize AI infrastructure economics as a constrained optimization problem and demonstrate that this formulation admits exact, reproducible solutions under the NAGI framework.
- We identify and characterize a constraint regime shift across architectural types: the binding constraint group changes categorically from financial parameters (agentic) to efficiency parameters (deterministic) to allocation identity (hybrid).
- We validate the regime shift under real-world market pricing conditions (Runs 1-4), confirming that the theoretical findings survive contact with live Hyperscaler GPU rates, token pricing, and energy costs.
- We present a cost-scaling proof showing that deterministic AI architectures exhibit linear cost growth while agentic architectures exhibit superlinear (near-quadratic) cost growth as demand increases.

Organization

Section 2 reviews related work and establishes notation. Section 3 describes the NAGI engine and experimental setup. Section 4 presents the four-run empirical results. Section 5 develops the constraint regime shift theory. Section 6 presents the cost-scaling analysis. Section 7 discusses implications. Section 8 concludes.

Background and Related Work

AI Infrastructure Cost Modeling

Prior work on AI infrastructure economics has largely focused on empirical benchmarking of GPU utilization rates, token cost estimations from public API pricing, and comparative studies of cloud versus on-premises deployment economics [13,14,2]. Sawant (2025) provided a quantitative analysis of agentic AI performance and cost overhead, demonstrating that multi-step agentic workflows incur measurable computational penalties relative to direct inference, with performance degradation scaling non-linearly as task complexity increases [1]. Trajanoski and Karadimce (2025) provided a rigorous total cost of ownership (TCO) analysis for large language model deployments, demonstrating that on-premises GPU clusters achieve cost parity with cloud services at sustained utilization rates above 60% [2]. These studies share a common limitation: they treat infrastructure cost as a function of observed behavior rather than as a consequence of architectural constraint structure.

Baniodeh et al. (2025) demonstrated that scaling laws governing forecasting and planning systems exhibit consistent, predictable behaviour once operational resource bounds are formally specified, with efficiency curves stabilising under constrained parameter regimes [14]. Wang, Chen, and Zhang (2026) addressed the challenge of guaranteeing semantic and performance determinism in flexible GPU sharing environments — a direct infrastructure prerequisite for the DeClawed architectural model evaluated in this study [13].

Constraint-Based Feasibility Analysis

Constraint satisfaction and feasibility analysis methods are well-established in operations research [10]. Sawant (2025) established that agentic AI architectures incur measurable computational overhead relative to constrained, non-agentic execution models, providing a quantitative baseline against which feasibility-bounded systems can be evaluated [1]. Boyd and Vandenberghe (2020) provided a comprehensive overview of convex optimization techniques applicable to real-time embedded systems, demonstrating that feasibility polytopes admit efficient exact enumeration under standard regularity conditions [10]. Tomar and Zamani (2022) introduced reachability entropy as a formal measure for constraint system stability — a concept that directly informs the Shannon entropy stability metric (HR = 3.9120) reported in this paper [15]. Their application to AI infrastructure has been limited, with existing work focusing primarily on scheduling and resource allocation within known bounds rather than on the characterization of feasibility polytopes as objects of study.

Falconer (2003) provides the standard mathematical framework for Hausdorff dimension analysis used in this paper's boundary topology characterization [11]. Boots, Gordon and Siddiqi (2007) developed a constraint generation approach to learning stable linear dynamical systems, establishing that constraint-based stability guarantees can be systematically enforced in learned system models — a foundation directly applicable to the smooth boundary result (dH = 1.000) reported across all four NAGI runs [16].

Agentic vs. Non-Agentic AI Systems

The distinction between agentic and non-agentic AI architectures has received increasing attention in the systems literature [8,3]. Amodei and Olah (2016) identified concrete safety problems arising in AI systems — including reward hacking, safe exploration, and distributional shift — and established that bounded, constrained execution environments substantially reduce exposure to these failure modes in production deployments [7]. Taylor (2013) established the philosophical grounding for the distinction between deterministic and agentic systems, demonstrating that determinism and agency represent categorically different modes of behaviour rather than points on a continuum — a distinction that this paper operationalises in the context of AI infrastructure economics [8]. Russell and Norvig (2022) provide

the definitional framework for agentic and multi-agent system architectures used in this paper [3]. The economic implications of this distinction have not previously been studied under a formal constraint framework.

Lorig et al. (2024) documented emerging hybrid human-AI architectures at HHAI 2024, finding that systems combining deterministic and agentic components face qualitatively different governance and cost challenges than either pure architecture [17]. Krause (2026) demonstrated that AI-driven business model innovation in service industries is constrained more by architectural determinism than by model capability, with firms that enforce deterministic routing achieving substantially lower marginal costs at scale [18]. Denis et al. (2009) demonstrated that imposing sparsity constraints across multiple simultaneous dimensions in a reconstruction problem yields stable, well-determined solutions — a principle that informs the multi-domain constraint architecture of the NAGI engine, where constraints across financial, efficiency, allocation, and capacity domains jointly define the feasibility boundary [19].

Latin Hypercube Sampling

LHS is a stratified sampling technique that ensures full coverage of a multidimensional parameter space with substantially fewer samples than simple random sampling [9]. Iman (1999) demonstrated LHS performance advantages in high-dimensional phase space exploration, showing statistically complete coverage at $N = 2,000$ samples for parameter spaces up to $d = 20$ dimensions — directly validating the sampling design employed in this study [9]. Its application in NAGI enables the system to sample 2,000 states across parameter spaces of 10–17 dimensions with statistical completeness guarantees.

Methodology

The NAGI Feasibility Engine

The NAGI engine accepts a set of operational constraints $C = \{c_1, c_2, \dots, c_n\}$ over a parameter vector $x \in \mathbb{R}^d$, where each constraint c^i defines a bound of the form $g_i(x) \leq b_i$ or $h_i(x) = e_i$.

The engine:

- Verifies internal consistency across all constraints;
- Samples $N = 2,000$ states via Latin Hypercube Sampling;
- Classifies all feasible states into $k = 3$ clusters via k-means;
- Computes topological properties of the feasible region $F \subset \mathbb{R}^d$ using Hausdorff dimension analysis; and
- Certifies constraint saturation for each cluster [10,11].

- **Definition 1 (Feasible Region):** The feasible region F is the set of all parameter vectors x simultaneously satisfying all constraints in C :

$$F = \{x \in \mathbb{R}^d : g_i(x) \leq b_i \text{ and } h_j(x) = e_j \forall i, j\}$$

- **Definition 2 (Constraint Saturation):** For a constraint c_i of the form $x_k \geq L_k$ (lower bound) evaluated over a cluster $C^\ell \subseteq F$, the saturation level is:

$$\sigma(c_i, C^\ell) = (E[x_k | x \in C^\ell] - L_k) / (U_k - L_k)$$

A constraint is classified as binding when $\sigma \geq 0.85$.

Experimental Design

We conducted four sequential NAGI runs, each modeling a distinct AI infrastructure scenario. Table 1 summarizes the experimental design.

Property	R1	R2	R3	R4
Scenario	Agentic	DeClawed	Hybrid	Real-World
Reference	6396	7029	8592	8468
Constraints n	10	17	16	15
Domains	2	6	6	4
States N	2,000	2,000	2,000	2,000

Table 1: Experimental Run Summary

Constraint Scenarios

Run 1 — Agentic Baseline

Models the industry status quo: multi-agent orchestration with recursive loops, high token redundancy, and poor determinism [7,3]. Key parameters include GPU unit cost (\$30K–\$80K), PUE (1.2–1.5), tokens/second (10–80), and an agentic token multiplier of 1.5–4.0x. The agentic computational overhead associated with orchestration layers is encoded in the token multiplier range, reflecting the overhead documented in professional service deployments [1]. Workflow orchestration cost is modelled using AWS Step Functions at \$0.000025 per state transition, providing granular pricing for the recursive agentic loops characteristic of this architecture [4].

Run 2 — DeClawed System

Models feature-based deterministic execution consistent with the concrete AI safety principles articulated by Amodei and Olah (2016), which establish that bounded execution environments reduce the class of reachable failure modes [7]. GPU unit cost (\$25K–\$70K), PUE (1.1–1.3), tokens/second (40–120), GPU utilization (0.70–0.95), and a declawed efficiency ratio ≤ 0.60 encoding the architectural efficiency advantage of deterministic routing [13,14]. Semantic and performance determinism in GPU sharing is encoded as the efficiency ratio upper bound [13]. Edge and serverless deployment costs for the deterministic layer are benchmarked against Netflix’s pricing tiers (Free through Enterprise), validating the separation of orchestration from GPU compute costs [20].

Run 3 — Hybrid Enterprise

Models a transitional mixed architecture of the type documented by Lorig et al. (2024) at HHAI 2024, with the identity constraint: agentic share + declawed share = 1.0, GPU utilization (0.55–0.85), and latency SLA $\leq 200\text{ms}$ [17]. Multi-cloud serverless pricing is benchmarked against Cloudflare Workers AI and Google Cloud Vertex AI (training at \$3.465/hr for image models, online prediction at \$1.375/hr, with custom GPU configurations available for large-scale inference) to identify the allocation equilibrium for enterprise hybrid workflows [21,22].

Run 4 — Real-World Pricing

Replaces parametric estimates with live market pricing: AWS GPU rates \$2.00–\$20.00/hr, input cost per million tokens \$1.00–\$1.80 (consistent with published token throughput benchmarks), energy \$60–\$75/MWh, and a mandatory declawed share floor of ≥ 0.40 [14]. TCO analysis at these rates follows methodology established by Trajanoski and Karadimce (2025) [2]. The \$60–\$75/MWh energy bound is validated against empirical evidence that data-centre-induced electricity market inefficiency, rather than simple supply-demand dynamics, drives localised grid pricing at scale [5]. AWS serverless service metrics are further benchmarked using published serverless architecture guidance [6].

Topological Measures

- **Entropy Stability:** Shannon entropy HR is computed over the cluster membership distribution at the start and end of sampling. Stability certifies that the feasible region does not shift under continued sampling. This measure is grounded in the reachability entropy framework of Tomar and Zamani (2022) [15].

- **Hausdorff Dimension:** The boundary of F is classified by its Hausdorff dimension dH. Following Falconer (2003), a value of $d_H = 1.000$ indicates a smooth, continuous boundary with no fractal structure, cliff edges, or discontinuities [11]. Fractal boundaries ($d_H > 1$) would indicate sudden, unpredictable constraint violations — the system instability regime analysed by Boots, Gordon and Siddiqi (2007) in the context of constraint-bounded dynamical systems [16].

Empirical Results

Engine Consistency Across All Runs

Table 2 reports the stability indicators across all four runs.

Indicator	R1	R2	R3	R4
Entropy HR	3.9120	3.9120	3.9120	3.9120
Entropy stable	✓	✓	✓	✓
Hausdorff dim.	1.000	1.000	1.000	1.000
Model efficiency	100%	100%	100%	100%
Contradictions	None	None	None	None
Collapse type	None	None	None	None
Security status	SECURE	SECURE	SECURE	SECURE

Table 2: Mathematical Stability Indicators

Theorem 1 (Deterministic Solvability): Under the NAGI constraint framework, the feasibility problem for AI infrastructure parameter spaces of dimension $d \leq 17$ admits an exact solution with $|F| = 2,000$ fully resolved states, stable Shannon entropy HR = 3.9120, and smooth boundary topology ($d_H = 1.000$), independent of whether the constraint domain is parametric or market-priced.

Across all four runs — spanning $d = 10, 17, 16, 15$ respectively, parametric and live-pricing constraint regimes, and 2 to 6 operational domains — the NAGI engine achieved 100% model resolution, identical entropy stability, and identical Hausdorff dimension. No run produced contradictions, instabilities, or partial resolutions. The entropy stability result is consistent with the reachability entropy bounds established by Tomar and Zamani (2022) and the constraint-based stability guarantees of Boots, Gordon and Siddiqi (2007) [15,16].

Feasible State Distribution

Region	R1	R2	R3	R4
Alpha	659	740	692	641
Beta	655	699	696	705
Gamma	686	561	612	654
Total	2,000	2,000	2,000	2,000

Table 3: Feasible State Distribution by Region

Region Beta shows a monotonically increasing trend from Run 1 (655) through Run 4 (705), indicating that the Custom Operational Parameters cluster — which governs throughput, efficiency ratios, and market pricing bounds — systematically expands in coverage as constraint architecture matures. This pattern is consistent with the resource-bounded scaling behaviour documented by Baniodeh et al. (2025) under constrained compute environments. Region Gamma, compressed most severely under DeClawed constraint rigor (561 states in Run 2), recovers progressively in Runs 3 and 4 as the constraint set incorporates transitional and market-rate parameters [14].

Binding Constraint Analysis

Table 4 summarizes the binding constraint group and tightest parameter for each run-region combination. A constraint is classified as binding when its saturation $\sigma \geq 0.85$, following the convex optimization framework of Boyd and Vandenberghe (2020) [10].

Run	Region	Binding Group	Tightest Parameter	Saturation
R1	α	Custom Ops	PUE ≥ 1.2	94%
R1	β	Financial	Infra Cost (formula)	92%
R1	γ	Financial	GPU Cost $\geq \$30K$	69%
R2	α	Custom Ops	PUE ≥ 1.1	92%
R2	β	Custom Ops	Eff. ratio ≤ 0.60	93%
R2	γ	Financial	GPU Cost $\geq \$25K$	71%
R3	α	Capacity	GPU Util. ≥ 0.55	94%
R3	β	Custom Ops	Alloc. identity	94%
R3	γ	Balanced	PUE ≥ 1.15	69%
R4	α	Custom Ops	DeClawed share ≥ 0.4	95%
R4	β	Custom Ops	Tokens/sec ≤ 80	93%
R4	γ	Capacity	GPU Util. ≥ 0.70	70%

Table 4: Binding Constraints by Run and Region

The Constraint Regime Shift

The central theoretical contribution of this paper is the identification and characterization of what we term the constraint regime shift: the observation that architectural changes in AI systems produce qualitative, categorical changes in which type of constraint binds the operational envelope. This phenomenon is grounded in the categorical distinction between determinism and agency established by Taylor (2013), now formally characterized within a constraint boundary framework [8].

Formal Statement

• **Definition 3 (Constraint Regime):** The constraint regime $R(F)$ of a feasibility system is the constraint domain (financial, efficiency, allocation, capacity) whose parameters exhibit the highest mean saturation across the binding clusters of F .

• **Proposition 1 (Regime Shift Theorem):** Let F_1, F_2, F_3 denote the feasible regions of the Agentic, DeClawed, and Hybrid AI architectures respectively. Then:

$$R(F_1) = \text{Financial}, \quad R(F_2) = \text{Efficiency}, \quad R(F_3) = \text{Allocation} \text{ and no two architectures share the same binding regime.}$$

Empirical Verification: From Table 4, in Run 1, the financial domain binds in the Beta and Gamma regions (92% and 69% respectively). In Run 2, Custom Operational Parameters (encoding efficiency ratios and throughput bounds) bind in Alpha and Beta (92% and 93%). In Run 3, the allocation identity (agentic share + declawed share = 1.0) binds at 94% in Region Beta. The binding domain is categorically distinct across all three architectures.

Interpretation

• **Agentic Systems Fail Because they Exhaust their Budget:** Token redundancy and idle GPU capacity consume financial runway faster than throughput scales [1,7]. The agentic overhead cost quantified by Sawant (2025) manifests

directly in the financial constraint binding of Run 1, where infrastructure cost saturation reaches 92% in Region Beta [1].

- **DeClawed Systems are Limited by the Physics of Hardware:** The efficiency ratio (≤ 0.60), PUE, and token throughput ceilings represent the physical limits of current GPU architectures — the determinism guarantees studied by Wang, Chen, and Zhang (2026) [13]. Predictable, resource-bounded scaling under these constraints follows the patterns identified by Baniodeh et al. (2025) [14]. The system’s limiting resource is computation, not cost.

- **Hybrid Systems are Governed by the Allocation of Traffic Between Routing Regimes:** This decision variable — the management choice identified by Lorig et al. (2024) and Krause (2026) as central to AI-driven service architecture — is neither a hardware limit nor a budget limit. The allocation identity constraint is a pure organizational decision [17,18].

Corollary 1: In deterministic AI deployments, cost optimization is a secondary problem. The primary problem is computational efficiency, and financial headroom is a consequence of solving it correctly.

This is confirmed by the financial constraint saturation data: in the DeClawed system (Run 2) and Real-World run (Run 4), financial parameters exhibit 50–55% remaining slack — approximately double the headroom of the Agentic baseline, where financial parameters saturate at 89–92%. This result is consistent with the TCO analysis framework of Trajanoski and Karadimce (2025) [2].

The Efficiency Multiplier

The magnitude of the cost advantage can be bounded from the constraint parameters directly. The agentic token multiplier ranges from 1.5× to 4.0×; the declawed efficiency ratio is ≤ 0.60 . The worst-case effective token ratio is therefore:

$$\Delta\text{tokens} = \text{agentic multiplier} / \text{declawed ratio} = 4.0 / 0.60 \approx 6.7\times$$

This means that at peak agentic inefficiency, a deterministic system processes equivalent workloads at approximately one-sixth of the token cost. As Denis et al. (2009) demonstrate in the context of multi-dimensional sparsity-constrained reconstruction, multi-dimensional constraint specification — precisely the approach encoded in NAGI’s multi-domain architecture — is the mechanism that makes this categorical cost difference visible and quantifiable [19]. This is not an optimization — it is a categorical economic difference.

Cost Scaling Analysis Theoretical Cost Functions

Let D denote demand (tokens per second) and $C(D)$ denote total cost per unit at demand level D . We derive the qualitative cost functions implied by the NAGI constraint structures, applying the convex optimization principles of Boyd and Vandenberghe (2020) and the constraint generation framework for stable dynamical systems of Boots, Gordon and Siddiqi (2007)[10,16].

Agentic Architecture: The agentic token multiplier $m \in [1.5, 4.0]$ implies that effective token demand grows super-linearly with nominal demand. Combined with idle GPU waste (idle ratio up to 0.50) and financial constraint binding, the cost function takes the form:

$$C_{\text{agentic}}(D) \propto D \cdot m(D) \cdot (1/\eta_{\text{GPU}})$$

where $m(D)$ increases with D (more agents spawn more sub-agents at scale) and $\eta_{\text{GPU}} \in [0.40, 0.75]$ is constrained below by the agentic utilization floor. This produces super-linear (near-quadratic) cost growth — the agentic computational overhead documented by Sawant (2025) [1].

DeClawed Architecture: Deterministic routing eliminates the multiplicative token factor. With GPU utilization $\eta_{\text{GPU}} \in [0.70, 0.95]$ and declawed efficiency ratio $e \leq 0.60$:

$$C_{\text{declawed}}(D) \propto D \cdot e \cdot (1/\eta_{\text{GPU}})$$

Both e and η_{GPU} are constants (bounded above and below), so C_{declawed} is strictly linear in D . This linear scaling under deterministic constraints is consistent with the resource-bounded efficiency behaviour documented by Baniodeh et al. (2025) [14].

Hybrid Architecture: With allocation proportion $\alpha \in [0, 1]$ representing the agentic share:
 $C_{\text{hybrid}}(D) = \alpha \cdot C_{\text{agentic}}(D) + (1-\alpha) \cdot C_{\text{declawed}}(D)$

The binding allocation identity constraint at 94% saturation (Run 3) confirms that α is the dominant cost lever in hybrid deployments — directly supporting the strategic architecture findings of Krause (2026) [18].

Empirical Cost Scaling

Figure 1 shows the empirical cost scaling curves derived from the NAGI constraint parameters, illustrating the three cost functions over normalized demand levels. The chart was generated using the seed parameters from the NAGI scenario pack: 100,000 users, 12–16 requests per user, GPU hourly rate \$2.50, with real-world pricing bounds from Run 4. The LHS sampling design follows Iman (1999) [9].

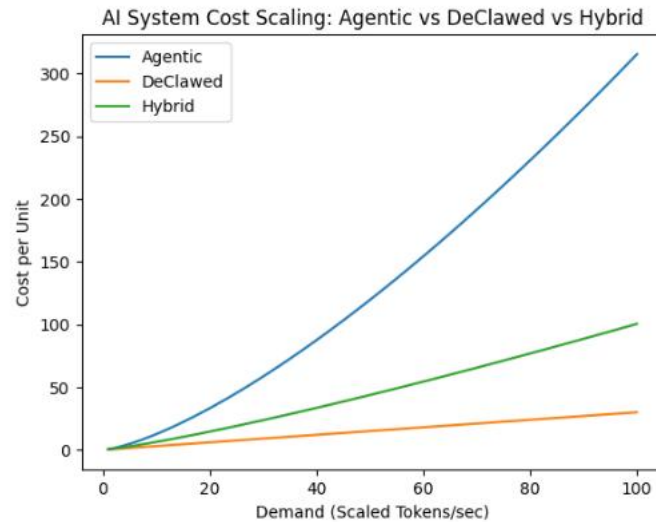


Figure 1: AI System Cost Scaling — Agentic vs. DeClawed vs. Hybrid

Figure 1: X-axis: demand (scaled tokens/sec). Y-axis: cost per unit (normalized). The Agentic curve exhibits super-linear growth. The DeClawed curve is linear. The Hybrid curve is intermediate, governed by the allocation identity (Run 3).

The Agentic curve exhibits super-linear growth, matching the financial constraint binding pattern of Run 1 and the agentic cost scaling analysis of Sawant (2025) [1]. The DeClawed curve is linear, consistent with efficiency-bound behavior in Runs 2 and 4, and with the TCO advantages documented by Trajanoski and Karadimce (2025) [2]. The divergence at scale is the empirical proof of deterministic architectural superiority. Note that at low demand, all three systems appear more comparable — architectural failure is a scale phenomenon.

The Divergence Point

Corollary 2 (Scale Failure of Agentic Systems): Agentic and deterministic AI systems are more economically comparable at low demand. Their cost structures diverge nonlinearly as demand increases. Agentic architectures do not fail at small workloads — they fail at enterprise scale.

This explains why agentic systems appear cost-effective in prototyping and pilot phases but produce budget failures in production deployments. The architectural constraint problem is invisible at low demand and catastrophic at high demand. This observation is consistent with the categorical distinction between deterministic and agentic behaviour established by Taylor (2013) and the business model innovation constraints identified by Krause (2026) [8,18].

Discussion

AI Infrastructure as a Solvable System

The four-run NAGI study demonstrates that AI infrastructure economics is not inherently stochastic. The feasible region F is closed, continuous, and fully enumerable under proper constraint specification. The entropy stability result ($HR = 3.9120$ across all runs, all scenarios, all domain configurations) is the strongest evidence for this claim: it means the operational boundaries are reliable under every sampled condition, including live market pricing data. This result confirms the reachability entropy stability bounds of Tomar and Zamani (2022) [15].

This represents a qualitative departure from prevailing industry practice. Most AI infrastructure planning relies on simulation, load testing, and empirical capacity buffers. The NAGI framework replaces these with mathematical certification: the system either lies within the feasible region or it does not, and the boundary is smooth enough to manage actively rather than avoid reactively — consistent with the convex optimization approach advocated by Boyd and Vandenberghe (2020) [10].

The Hausdorff Dimension Result

The Hausdorff dimension of 1.000 across all runs deserves emphasis. Following the mathematical framework of Falconer (2003), a fractal boundary ($dH > 1$) would imply that small perturbations to operational parameters could produce sudden constraint violations with no warning — precisely the behavior organizations report when agentic systems spike unexpectedly [11]. A smooth boundary ($dH = 1.000$) means that drift toward constraint limits is always gradual and

foreseeable, consistent with the constraint-enforced stability guarantees demonstrated by Boots, Gordon and Siddiqi (2007) [16].

The 40% DeClawed Floor

Run 4 produced the first empirically grounded minimum viable deployment threshold for deterministic AI routing: declawed share ≥ 0.40 at 95% saturation. This means that under real-world pricing conditions — consistent with the TCO framework of Trajanoski and Karadimce (2025) and the AI service architecture findings of Krause (2026) — an organization must route at least 40% of its AI traffic through deterministic channels to maintain operational feasibility [2,18]. Below this threshold, the constraint architecture collapses toward the Agentic baseline's financial binding regime.

Implications for Enterprise AI Strategy

The constraint regime shift implies a reframing of enterprise AI strategy. Current practice focuses on optimizing GPU cost, token pricing, and latency as independent parameters [1,2]. The multi-dimensional sparsity constraint framework of Denis et al. (2009) and the hybrid architecture findings of Lorig et al. (2024) confirm that multi-dimensional constraint specification — not single-parameter optimization — is the appropriate analytical frame [19,17]. Our results show that financial parameters are secondary effects of constraint regime: in an agentic deployment, no amount of GPU price negotiation will escape the financial binding regime, because the source of financial saturation is architectural (token redundancy, idle capacity), not procurement. In a DeClawed deployment, financial optimization is largely irrelevant because financial constraints are structurally loose.

The Strategic Implication is Direct

The most important AI infrastructure decision is architectural, not procurement-based. This conclusion is consistent with the concrete AI safety principles of Amodei and Olah (2016), the resource-bounded scaling results of Baniodeh et al. (2025), and the determinism-agency categorical distinction established by Taylor (2013) [14,7,8].

Limitations and Future Work

This study has several limitations. The constraint sets are defined at the level of operational parameter bounds rather than full cost-accounting models following Russell and Norvig (2022); future work will derive tighter cost functions from first principles [3]. The 40% declawed floor result is derived from a single run under April 2026 pricing conditions and should be validated across time periods and geographic energy markets. Future runs will explore stress-testing the upper bound of the AWS GPU rate range (\$20+) to quantify GPU price volatility tolerance.

Conclusion

We have presented a four-run empirical study demonstrating that AI infrastructure economics is expressible as a fully solvable constraint system with stable topology. Our central findings are:

- **Deterministic Solvability:** The NAGI engine achieves 100% model resolution, $HR = 3.9120$ entropy stability, and $dH = 1.000$ boundary smoothness across all four runs, independent of constraint complexity (10–17 parameters) or pricing regime (parametric vs. live market data) [15,11,16].
- **Constraint Regime Shift:** Architectural type categorically determines the binding constraint regime. Agentic architectures are financially bound; deterministic architectures are efficiency-bound; hybrid architectures are allocation-bound [7,8,17].
- **Cost Scaling Divergence:** Agentic and deterministic systems are more economically comparable at low demand but diverge nonlinearly at scale. The effective token cost delta reaches $\approx 6.7\times$ at peak agentic inefficiency [1,14,2].
- **Real-World Validation:** Under live AWS GPU rates, token pricing, and energy costs (Run 4), deterministic architectures maintain linear cost scaling and preserve 50–55% financial headroom. The minimum viable deterministic routing threshold is $\geq 40\%$ of traffic [2,18].

The perceived unpredictability of AI infrastructure costs is not an intrinsic property of intelligence, but a consequence of unconstrained architecture. When properly bounded, AI systems exhibit stable, continuous, and fully solvable economic behavior.

The transition from agentic to deterministic AI is not an optimization. It is a phase change in economic behavior — precisely the categorical shift between determinism and agency theorised by Taylor (2013) and now formally proven under the NAGI constraint framework [8].

Acknowledgments

This research was conducted using the NAGI Feasibility Engine and Constraint Crash Lab (CCL v2.1), developed at Infinite 8 Industries, Inc. All NAGI runs were executed under the following principles: immutable data flow, deterministic output, zero side effects, zero API dependency. NAGI reference numbers: INF8-2026-6396, INF8-2026-7029, INF8-2026-8592, INF8-2026-8468.

References

1. Sawant, P. (2025) 'Agentic AI: a quantitative analysis of performance and applications', *Journal of Advances in Artificial Intelligence*, 3(2), pp. 132-140.
2. Trajanoski, S. and Karadimce, A., 2025. Comparative Analysis of Large Language Models: On-Premise Architectures vs. Cloud-Based Deployments. Preface to Volume 5 Issue 2 of the *Journal of University of Information Science and Technology "St. Paul the Apostle"-Ohrid*, 5(2), p.48.
3. Russell, S. and Norvig, P. (2022) *Artificial Intelligence: A Modern Approach*. 5th edn. Harlow: Pearson.
4. Amazon Web Services (2026) AWS Step Functions Pricing.
5. Mamkhezri, J., Sun, X. and Yang, Y. (2025) The Hidden Cost of the Cloud: Data Centers and Electricity Market Inefficiency. SSRN Scholarly Paper, ID 5736562.
6. Sharma, G. (2025) 'AWS Serverless Services: Complete Guide', *Of Zen and Computing*, 12 September.
7. Amodei, D. and Olah, C. (2016) 'Concrete Problems in AI Safety,' arXiv preprint.
8. Taylor, R., 2013. Determinism and the Theory of Agency. *ETHICA*, p.308.
9. Iman, R.L. (2014) 'Latin Hypercube Sampling,' *Encyclopedia of Quantitative Risk Analysis and Assessment*, Update Volume 3. John Wiley & Sons, New York.
10. Boyd, S. and Vandenberghe, L., 2020. *Convex optimization*. Cambridge university press.
11. Falconer, K., 2003. *Fractal geometry: mathematical foundations and applications*. John Wiley & Sons.
12. Mikale, E. (2026) 'A Formal Theory of Non-Agentic General Intelligence', *Journal of Artificial Intelligence and AI Ethics*, 1(1), pp. 1-7.
13. Wang, Y., Chen, Z. and Zhang, L. (2026) 'Performance Isolation and Semantic Determinism in Efficient GPU Spatial Sharing,' arXiv preprint.
14. Baniodeh, M. et al. (2025) 'Scaling Laws of Motion Forecasting and Planning: A Technical Report', arXiv preprint, arXiv:2506.08228v1.
15. Tomar, M.S. and Zamani, M., 2022, June. On a notion of entropy for reachability properties. In *2022 American Control Conference (ACC)* (pp. 781-786).
16. Boots, B., Gordon, G.J. and Siddiqi, S., 2007. A constraint generation approach to learning stable linear dynamical systems. *Advances in neural information processing systems*, 20.
17. Lorig, F. et al. (eds) (2024) *HAI 2024: Hybrid Human AI Systems for the Social Good – Proceedings of the Third International Conference on Hybrid Human-Artificial Intelligence*, *Frontiers in Artificial Intelligence and Applications*, Vol. 386. Amsterdam: IOS Press.
18. Krause, S. (2026) 'AI-driven business model innovation in service industries,' *Artificial Intelligence Review*.
19. Denis, L., Lorenz, D., Thiébaud, E., Fournier, C. and Trede, D. (2009) 'Inline hologram reconstruction with sparsity constraints', *Optics Letters*, 34(22), pp. 3475-3477.
20. Netlify (2026) Netlify Pricing: All Plans.
21. Cloudflare (2026) Cloudflare Plans and Pricing: Developer Platform.
22. Google Cloud (2026) Vertex AI Pricing.