


# When Benchmarks Saturate: Evaluation-Deployment Structural Divergence in Large-Scale AI Systems

Gregor Herbert Wegener 

Friedrichstrasse 4, 10969 Berlin, Germany; gregor.wegener@gmail.com; Tel.: +49 179 2544522

## Abstract

Benchmark scores across frontier AI models are converging on major evaluation suites, while deployment behavior is becoming increasingly shaped by serving conditions, runtime coordination, tool orchestration, and context management. This shift changes the practical meaning of model comparison. As score differentials narrow, capability measurements alone provide less explanatory power for production decisions, and execution geometry becomes more relevant for understanding how systems behave under real operating conditions. Evaluation can be understood as a projection of model behavior onto a controlled context defined by task format, inference settings, and bounded runtime assumptions. Deployment contexts differ structurally from these conditions through heterogeneous infrastructure, dynamic serving policies, recursive workflows, and load dependent execution paths. The resulting divergence is not well captured by classical benchmark methodology, because benchmarks measure task performance under fixed conditions, while production behavior emerges from the interaction between model, serving stack, and infrastructure. This paper introduces a constructive diagnostic perspective for analysing evaluation deployment divergence in large scale AI systems. It maps this perspective to the SORT application catalog through ai.47, ai.16, and ai.02, and frames the contribution as an extension of current evaluation practice that can support greater transparency, comparability, and deployment robustness in hyperscale inference environments.

**Keywords:** benchmark saturation; evaluation deployment divergence; evaluation context projection instability; structural drift; inference geometry; runtime coherence; benchmark integrity; structural diagnostics; hyperscale AI infrastructure; deployment robustness

---

## 1. Executive Summary

Benchmark scores across frontier AI models are converging on major evaluation suites, while frontier releases continue to advance reasoning, coding, and multimodal performance at high speed [1, 14]. As these score differentials narrow, their value as a primary decision variable for production deployment declines. Small benchmark advantages can remain informative at the task level, but they increasingly provide limited visibility into how systems will behave once they are embedded in real serving environments.

This convergence should not be interpreted as evidence that leading models have become operationally interchangeable. Rather, it indicates that the evaluation geometry through which models are compared is approaching saturation. Under these conditions, the decisive variable for deployment begins to shift from benchmarked capability alone toward execution geometry, namely the way model behavior is shaped by serving conditions, runtime coordination, tool orchestration, and context management in production settings [2].

Recent model releases illustrate this transition clearly. GPT-5.3 Codex is positioned around extended coding workflows with tool use and multi-step execution [4]. Claude Opus 4.6 emphasizes prolonged agentic workflows, computer-use interfaces, and large-context execution [5]. Gemini 3.1

Pro introduces dynamic inference depth allocation across multiple reasoning tiers [6]. In each case, practical differentiation is shaped not only by model capability in isolation, but by the structure through which that capability is expressed at runtime.

The central structural observation of this paper is that evaluation contexts and deployment contexts project onto different regions of a model's behavioural space. A model that performs strongly under controlled benchmark conditions may exhibit materially different behaviour under production-scale serving, heterogeneous hardware, runtime coordination, and recursive workflow conditions. This divergence is not best understood as random variation. It emerges from the structural difference between the context in which behaviour is measured and the context in which behaviour is operationalized. In heterogeneous inference systems, the serving stack itself functions as a transformation layer that can reshape execution geometry independently of model weights [24].

Classical evaluation methodology provides only partial visibility into this shift. Benchmarks measure capability under bounded and comparatively static conditions. Logs capture outputs and events, but not the full execution topology that gives rise to those outputs. Production monitoring tracks indicators such as utilization, latency, and throughput, yet these remain layer-specific views that do not directly represent the coupled interaction surface across accelerator, network, virtualization, memory, and scheduling layers where structurally relevant behaviour becomes visible [24].

This creates a growing gap between what evaluation confirms and what deployment must actually sustain. As systems become more dependent on serving architecture, infrastructure heterogeneity, and multi-step orchestration, the evaluation-deployment boundary becomes less of a clean handoff and more of a structural coupling surface. Under benchmark saturation, this surface gains strategic relevance because it increasingly conditions whether nominal model capability is translated into robust, efficient, and predictable production behaviour.

The purpose of this paper is to introduce a constructive diagnostic perspective for analysing evaluation-deployment structural divergence and to relate this perspective to the SORT application catalog. In particular, the discussion is anchored in ai.47, ai.16, and ai.02 as complementary views on projection instability, benchmark drift, and structural drift across AI workloads. The objective is not to displace existing benchmark practice, but to extend it with a systems-level vocabulary that can support greater transparency, comparability, and deployment robustness in hyperscale AI environments.

## 2. The Economic Shift: From Capability to Geometry

### 2.1. Benchmark Convergence and the Saturation Threshold

Frontier model development is continuing at a high pace, and leading systems are increasingly being compared within relatively narrow performance bands across widely used evaluation suites [1,5,6]. This does not imply that frontier models have become equivalent. It indicates instead that score differentials on established benchmarks are becoming less decisive as a standalone basis for deployment choice. Under these conditions, benchmark outcomes remain useful for capability screening, but they provide less explanatory power for how systems will perform once they are embedded in production environments.

At the same time, competitive differentiation is shifting away from raw model quality in isolation and toward economically relevant execution characteristics. Inference cost trajectories, responsiveness under load, integration depth, and support for longer-running workflows are becoming more visible determinants of adoption than isolated score advantages alone [2]. As reasoning depth, coding assistance, and multimodal competence become more broadly available at the frontier, the operational meaning of performance is increasingly defined by how capability is expressed through runtime conditions rather than by benchmark rank alone.

This shift becomes structurally significant once benchmark compression reaches a point at which differences remain measurable but are no longer strongly decision-relevant on their own. In that regime, organisations are less able to select models primarily through capability comparison and more likely to differentiate on deployment robustness, serving compatibility, and the stability of behaviour

across real operating conditions. The threshold is therefore not purely numerical. It is reached when the marginal decision value of benchmark superiority declines relative to the operational value of execution geometry.

## 2.2. *The Inference Infrastructure Shift*

The broader economic environment of AI deployment is also changing. Across the market, more attention is moving from training scale alone toward the cost, efficiency, and controllability of inference at production scale [2,3]. At the same time, datacenter expansion, power demand, and infrastructure planning are rising in strategic importance as AI adoption deepens across enterprise and platform environments [17,18]. These developments do not reduce the significance of training. They increase the importance of the serving layer through which trained capability is operationalized.

Industry moves in early 2026 reinforce this interpretation. The OpenAI–Amazon partnership places explicit emphasis on stateful runtime environments and enterprise-oriented distribution pathways [8]. NVIDIA’s inference-focused accelerator direction and TensorRT-LLM AutoDeploy similarly point toward a deployment regime in which runtime automation, serving optimisation, and heterogeneous execution stacks become central engineering concerns [9,10]. In such a regime, production performance depends not only on the model artifact, but on the infrastructural conditions that shape its execution.

The strategic implication is that model capability is becoming a necessary but not sufficient condition for deployment success. As inference infrastructure grows more heterogeneous and operational requirements grow more demanding, infrastructure architecture, serving stack design, and runtime coordination emerge as co-determinants of realised effectiveness. Under benchmark saturation, this makes the infrastructure layer more visible as a differentiating surface in its own right.

## 2.3. *From Model Selection to Execution Geometry*

The transition from capability-driven comparison to geometry-sensitive deployment has practical consequences for how frontier systems are interpreted. Architectures such as Qwen 3.5, which combine very large total parameter counts with a smaller active parameter footprint per forward pass, illustrate how inference efficiency and execution structure can become strategically relevant when workflows extend beyond single-step prompting into sustained interaction patterns [7]. In such settings, model choice is increasingly inseparable from the runtime conditions under which the model is expected to operate.

A similar point is visible in systems that expose dynamic inference behaviour directly. Gemini 3.1 Pro’s multi-tier reasoning allocation shows that model performance is no longer best understood as a single static property [6]. The same model can generate different execution profiles depending on runtime configuration, latency budget, or reasoning depth. These distinctions matter in deployment, yet they are only partially visible in single-context benchmark measurements. What appears as a single model in evaluation may function as multiple operational regimes in production.

This becomes even more pronounced in agentic settings. Industry reporting suggests that many agent-oriented initiatives face difficulty in reaching stable production use, while recursive workflows can expand token usage and orchestration overhead substantially relative to simpler interaction patterns [19,20]. These observations point toward an important architectural reality. In longer-running systems, realized performance is shaped not only by model capability, but by whether execution remains efficient, coordinated, and predictable as workflow complexity grows.

For this reason, model selection is increasingly becoming a question of execution geometry rather than capability alone. Research and deployment practice are moving toward tighter co-design between model, serving stack, and hardware [10]. At the same time, larger cloud and chip partnerships introduce new dependency structures across runtime environments and distribution channels [8]. The evaluation-deployment boundary is therefore becoming less of a simple transition step and more of a structurally important interface through which production behaviour is shaped.

### 3. The Hidden Structural Effect: Evaluation-Deployment Divergence

#### 3.1. Evaluation as Context-Dependent Projection

Evaluation is not a neutral reading of model capability in the abstract. It is a projection of model behaviour onto a specific evaluation context defined by dataset composition, prompt structure, inference parameters, hardware assumptions, batch conditions, and bounded runtime settings [14]. This framing does not reduce the value of benchmarking. It clarifies that benchmark results are always produced within a particular measurement geometry.

Under benchmark conditions, that geometry is typically designed to be controlled, reproducible, and comparable across models. This makes evaluation highly useful for task-level differentiation. At the same time, it means that the measured behaviour corresponds to a selected region of the model's broader behavioural space. Production systems project the same model through different conditions, including serving policies, orchestration logic, tool use, context persistence, infrastructure heterogeneity, and load-sensitive execution paths. The evaluated model and the deployed model therefore operate through structurally different contexts even when the underlying weights remain unchanged.

The divergence between these projections is best understood as a structural effect rather than as random variance. It arises because evaluation and deployment do not expose the model to the same execution conditions or the same coordination surfaces. This is the central concern captured in the SORT catalog by ai.47, *Evaluation Context Projection Instability*, which addresses behaviour divergence between evaluation and deployment contexts as a distinct structural phenomenon.

#### 3.2. Mechanisms of Structural Divergence

Several current optimisation patterns make this divergence more visible. Speculative decoding and related serving-level optimisations can reshape inference execution geometry without changing the model itself. SPEED-Bench is notable precisely because it addresses speculative decoding under more realistic serving regimes, thereby recognising that optimisation techniques introduce data-dependent and context-sensitive performance profiles that are not fully represented in classical evaluation setups [11].

A similar shift appears in agentic execution. Systems such as GPT-5.3 Codex and Claude Opus 4.6 extend execution from bounded single-pass inference into multi-step workflows involving tool use, retry behaviour, and context management over longer horizons [4,5]. In these settings, execution quality depends not only on model competence in isolation, but also on whether recursive workflows remain coordinated as interaction depth increases [21]. Benchmarks that evaluate single-turn capability can therefore understate the extent to which runtime structure shapes realized behaviour in production.

Further divergence is introduced by serving-time decisions such as context truncation, KV-cache management, and dynamic batching. These are control-layer mechanisms that improve efficiency and throughput, but they also modify the effective execution conditions under which the model operates. When independently acting coordination layers interact without complete global visibility, coherence can become fragmented in ways that are not legible through fixed-context benchmark methodology alone [22]. The result is not necessarily degradation in the ordinary sense. It is a change in the behavioural geometry through which the model's outputs are produced.

The recent emergence of evaluation efforts such as ISO-Bench and FeatureBench is therefore significant [12,13]. Both move toward more realistic workload conditions, coding-agent settings, and drift-aware task construction. This is a constructive development. It increases evaluation resolution at the serving and workflow layer. At the same time, these benchmarks should be understood as partial approximations of deployment structure rather than full visibility into evaluation-deployment divergence.

### 3.3. *Serving Stack as Behavioural Modifier*

The serving stack is often treated operationally as the delivery layer through which model capability is exposed. Structurally, however, it functions as a transformation layer. Batch scheduling, memory management, request routing, accelerator assignment, and virtualization boundaries all shape the path through which requests are executed. In heterogeneous inference systems composed of mixed accelerator classes, disaggregated serving paths, and virtualized runtime layers, these transformations can produce cross-layer incoherence even when individual components remain locally well-behaved [24].

This means that the same model can produce different effective execution geometries under different serving configurations. Behaviour becomes accelerator-dependent, region-dependent, load-dependent, and policy-dependent. Such variation should not be interpreted as anomalous in itself. It is a normal structural property of heterogeneous inference environments in which runtime composition matters as much as component quality. What changes is not necessarily the model artifact, but the execution surface through which that artifact becomes operational.

Current industry developments increasingly reflect this reality. The OpenAI–Amazon partnership places emphasis on stateful runtime environments and enterprise-oriented deployment pathways [8]. NVIDIA’s AutoDeploy work similarly points toward automated deployment optimisation across heterogeneous stacks [10]. Together, these developments suggest growing recognition that serving infrastructure is not merely an efficiency layer. It is part of the behavioural interface of the deployed system. Evaluation practice has begun moving in this direction, but it has not yet fully integrated the serving stack as a first-order analytical variable.

### 3.4. *The Saturation Amplifier*

The effect of evaluation–deployment divergence becomes more important as benchmark scores converge. When benchmark gaps are large, capability differences dominate model selection and can outweigh moderate differences in deployment structure. In that regime, structural divergence remains present, but it is often operationally secondary. As score differentials narrow, this balance changes. The relative importance of execution geometry increases because capability alone offers less separation between candidate systems.

This creates an amplification effect. The closer benchmark outcomes cluster, the more production behaviour depends on the structural conditions through which capability is expressed. Benchmark saturation therefore does not reduce the relevance of evaluation. It increases the importance of alignment between evaluation context and deployment context. The strategic question is no longer only which model scores highest under controlled conditions, but which model-context combination remains most stable, efficient, and legible under real operating conditions.

For deployment teams, the practical consequence is straightforward. When decisions are made primarily on the basis of saturated benchmark scores, the selected variable may be losing discriminatory power just as execution geometry is gaining it. A more structurally informed view does not replace benchmarking. It extends benchmarking by clarifying when task-level comparison is sufficient and when context-sensitive deployment behaviour becomes the more relevant basis for selection.

## 4. Why Benchmarks Don’t Catch This

### 4.1. *Benchmarks Measure Capability, Not Projection Stability*

Classical benchmarks are designed to answer questions of task competence. In practical terms, they ask whether a model can solve a defined task under specified conditions. They are generally not designed to determine whether the same model remains structurally stable as execution context changes across serving conditions, orchestration logic, hardware configuration, or workflow depth [14, 15]. This is not a weakness in benchmark design. It reflects the fact that most benchmark methodologies are optimized for comparability, reproducibility, and task-level discrimination rather than for context-sensitive systems analysis.

The structural limitation arises because benchmarks project behaviour onto a bounded evaluation space under controlled assumptions. That space can be expanded, diversified, or made more realistic, but it remains an evaluation space. As a result, drift in dimensions that are activated only under deployment conditions may remain only partially visible or entirely unobserved. The core issue is therefore not simply incomplete benchmark coverage. It is the more general limitation of point-in-time measurement when applied to systems whose realised behaviour depends on changing execution structure.

Under these conditions, higher benchmark performance can coexist with lower deployment stability if the evaluation context is only weakly representative of the deployment context. Benchmark saturation makes this more consequential. When score gaps narrow, differences that appear small in evaluation may still correspond to materially different behaviour once models are embedded in production systems with heterogeneous infrastructure and runtime adaptation. The benchmark result remains real, but its decision value becomes more dependent on the structural relation between evaluation and deployment.

#### *4.2. New Benchmarks Improve Resolution, Not Structural Visibility*

Recent benchmark development is already moving in a constructive direction. SPEED-Bench addresses speculative decoding under diverse workloads [11]. ISO-Bench targets real inference workloads and coding-agent optimisation [12]. FeatureBench evaluates end-to-end feature development with explicit attention to drift, leakage, and task realism [13]. WebArena provides more realistic autonomous agent environments than conventional static task suites [16]. Together, these efforts represent an important increase in evaluation realism.

What these advances primarily improve, however, is evaluation resolution rather than full structural visibility. They bring benchmarking closer to the operational layer by incorporating more realistic workloads, richer task structure, and more demanding execution settings. This is highly valuable. At the same time, a benchmark can show that performance changes under a more realistic condition without fully revealing how the execution structure changed to produce that result. It can detect behavioural difference without necessarily exposing the topology of the transformation that generated that difference.

For this reason, better benchmarks should be understood as refinements of task-level evaluation rather than as substitutes for structural observability. They improve the quality of comparison and reduce some forms of abstraction error. They do not by themselves make the full evaluation-deployment projection legible. That additional visibility requires a diagnostic layer capable of relating measured outcomes to the execution conditions through which those outcomes were produced.

#### *4.3. Dashboards Track Performance, Not Topology*

Production monitoring systems provide indispensable operational visibility, but they do so through layer-specific metrics such as utilization, latency, throughput, queue depth, and error rate. These measurements are necessary for operating large-scale systems effectively. Their limitation is different: they generally describe local or aggregated performance states without directly representing the coupled interaction surface across accelerator runtime, memory path, virtualization boundary, network topology, and scheduling logic on which structurally relevant behaviour emerges [24].

As a result, systems can appear healthy in dashboards while already operating in a structurally altered regime. Service indicators may remain acceptable even as configuration-sensitive drift accumulates across releases, infrastructure changes, routing policies, or serving-layer optimisations. In such cases, benchmark-verified performance and operational performance are not simply separated by time. They are separated by a transformation layer that remains only partially visible to standard telemetry. The issue is not insufficient instrumentation in the ordinary sense. It is that topology-sensitive behaviour is not directly represented by the metrics most commonly observed.

The missing dimension is therefore not more metrics alone, but a clearer structural mapping between what benchmarks measure and what deployment conditions actually produce. That mapping

is the evaluation-deployment projection. Its stability is the central concern of ai.47, *Evaluation Context Projection Instability*. A structurally informed perspective does not replace dashboards or benchmarking. It complements both by making the relation between measured capability and realised execution behaviour more explicit.

## 5. The Structural Diagnostic Gap

### 5.1. What Current Evaluation Does Not Capture

Current evaluation methodology is well suited for measuring whether a model produces acceptable outcomes under specified conditions. What it captures less directly is whether the same query follows different execution trajectories once deployment conditions vary across hardware class, routing logic, load state, batch policy, or serving configuration. In other words, evaluation is primarily outcome-oriented. It is generally not designed to expose the structural path through which those outcomes are produced.

This limitation becomes more relevant as execution environments grow more configurable and heterogeneous. A change in deployment configuration can preserve benchmark-visible performance while still altering the underlying execution geometry in ways that affect reproducibility, efficiency, or stability in production. Across repeated releases and infrastructure adjustments, such shifts can accumulate into meaningful behavioural divergence without appearing clearly in benchmark results, because the benchmark continues to validate task capability rather than structural continuity.

The engineering community is increasingly moving toward this observation. Work on speculative decoding realism, serving-sensitive benchmarking, and drift-aware benchmark construction suggests growing recognition that evaluation outcomes can depend materially on execution context rather than on model capability alone [11,13]. This is a constructive development. It improves awareness of the gap. What remains less developed is a diagnostic layer capable of systematically tracking the structural conditions under which that gap expands or contracts.

### 5.2. Context Divergence as a Structural Property

The divergence between evaluation and deployment contexts is best understood as a structural property of large-scale AI systems rather than as a temporary shortcoming of benchmark design. More comprehensive evaluation can improve realism, but it does not eliminate the fact that evaluation and deployment are organized around different kinds of context. One is built for reproducibility and comparability. The other is shaped by variability, infrastructure coupling, and runtime adaptation.

Evaluation contexts are intentionally bounded. They are constructed to make performance legible across models, tasks, and versions. Deployment contexts are structurally different. They are influenced by serving policies, orchestration decisions, accelerator heterogeneity, memory constraints, virtualization overhead, and changing workload conditions. These contexts therefore do not simply differ in detail. They differ in operating logic. The relevant question is not whether divergence can occur, but how it can be made more explicit, observable, and interpretable over time.

Seen in this way, context divergence is not something that disappears once enough benchmark tasks are added. It remains present because the measurement environment and the operating environment serve different purposes. This is why the problem cannot be fully resolved through benchmark expansion alone. What is additionally required is a way to represent the structural relation between the two contexts and to monitor whether that relation remains stable as systems evolve.

### 5.3. The Missing Layer

The unresolved gap is therefore a structural diagnostic layer for evaluation-deployment divergence. Such a layer would not ask only whether a model succeeds on a benchmark. It would ask whether the projection from evaluation context to deployment context remains sufficiently stable for benchmark results to retain decision value under changing serving conditions.

This kind of layer would complement existing evaluation and monitoring rather than replace them. Benchmarks would continue to provide task-level comparison. Telemetry would continue to provide operational visibility. The added contribution would be to connect these two views by making the context-dependent nature of realised performance more legible. In practical terms, this means improving visibility into how serving configuration, infrastructure state, and execution topology condition the translation from measured capability to deployed behaviour.

For hyperscalers and large-scale AI platform operators, this should be understood as a constructive extension of current practice. As benchmark saturation progresses and deployment environments become more structurally differentiated, one of the next useful improvements in evaluation methodology is likely to come from greater structural visibility into how benchmark results project onto production conditions. That visibility is the missing layer addressed in this paper.

## 6. SORT as a Structural Diagnostic Layer

### 6.1. Diagnostic Vocabulary

Projection Instability.

Projection instability denotes the structural condition in which a model's evaluated behaviour does not reliably predict its deployment behaviour because evaluation and deployment contexts project onto structurally different behavioural regions. This concept is central to ai.47, *Evaluation Context Projection Instability*, and provides the primary diagnostic lens for this paper. The Public Application Catalog defines ai.47 as the structural analysis of behaviour divergence between evaluation and deployment contexts, including patterns in which benchmark-visible performance does not transfer cleanly to production conditions. At the mechanism level, heterogeneous inference infrastructure provides one concrete pathway through which serving conditions reshape execution geometry without changing model weights [24].

Benchmark Drift.

Benchmark drift denotes the gradual loss of benchmark relevance as a system's structural behaviour evolves across releases, configurations, and infrastructure conditions without corresponding evolution in the benchmark suite. A system may continue to perform well under established benchmark conditions while its production-relevant behaviour changes in dimensions not directly represented in those evaluations. This concept maps to ai.16, *Benchmark Integrity and Drift Diagnostics*, which frames structural stability metrics as a complement to classical benchmarks for detecting drift across releases and configurations.

Structural Drift without Model Modification.

Structural drift without model modification denotes the reshaping of effective model behaviour through changes in serving stack, infrastructure configuration, runtime coordination, or deployment topology while the model weights remain unchanged. This maps to ai.02, *Structural Drift Diagnostics for AI Workloads*, which addresses structural drift across training and inference pipelines beyond metrics and telemetry. In practice, such drift can arise when cross-layer interactions between accelerator runtimes, network topology, virtualization boundaries, and scheduling logic alter the execution conditions through which the model is expressed [24].

Evaluation Context Coupling.

Evaluation context coupling denotes the degree to which benchmark results depend on specific evaluation conditions rather than reflecting more invariant model properties. High coupling means that comparatively small changes in prompt structure, runtime settings, or task framing can produce disproportionately large changes in measured performance. Under benchmark saturation, this becomes more relevant because the remaining score differentials may partly reflect evaluation-context artefacts rather than durable capability differences.

## 6.2. Application Mapping

**Table 1.** SORT Application Catalog Mapping for Evaluation-Deployment Divergence

ID	Cluster	Application	Structural Mapping
ai.47	C	Evaluation Context Projection Instability	Evaluation contexts project onto a specific region of the model’s behavioural space. When this projection is structurally unrepresentative of the deployment context, evaluation results lose predictive validity for production behaviour. At the mechanism level, serving infrastructure can reshape execution geometry through configuration-dependent runtime conditions [24].
ai.16	B	Benchmark Integrity and Drift Diagnostics	Structural stability metrics complement classical benchmarks by detecting drift across releases and configurations in dimensions not captured by standard benchmark suites. Temporal drift can accumulate silently and widen the gap between benchmark-verified performance and actual operational stability.
ai.02	A	Structural Drift Diagnostics for AI Workloads	Infrastructure-induced behavioural change without model modification is detectable through structural drift analysis across training and inference pipeline configurations. Physical coupling and serving-layer reconfiguration can contribute to this form of drift as execution conditions change over time [23,24].

## 7. Strategic Implications for Hyperscalers

### 7.1. Deployment Robustness as the New Differentiator

When benchmark scores converge, deployment decisions are shaped less by the question of which model scores marginally higher under controlled evaluation and more by which model infrastructure combination delivers the most stable, efficient, and predictable behaviour under production conditions. This elevates serving architecture, runtime coherence, and infrastructure model compatibility from secondary implementation considerations to primary differentiators in large scale inference environments.

For hyperscalers operating multi model and multi accelerator fleets, evaluation deployment alignment therefore becomes a fleet level operational capability rather than a narrow pre deployment checkpoint. In heterogeneous serving environments, runtime coherence across accelerator, network, virtualization, and scheduling layers conditions whether provisioned capacity is translated into effective capacity at the system level [24]. Under these conditions, deployment robustness functions as a practical expression of structural fit between model behaviour and execution environment.

### 7.2. Evaluation Methodology Requires Structural Extension

Current evaluation pipelines are highly effective at measuring capability under controlled conditions. As benchmark saturation progresses, however, their practical value increasingly depends on how clearly evaluation results can be related to production conditions. This suggests the need for a structural extension of evaluation methodology, namely a more systematic assessment of how benchmark outcomes project onto serving environments, orchestration patterns, and infrastructure specific execution paths.

Recent benchmark development already points in this direction. SPEED-Bench, ISO-Bench, and FeatureBench each introduce more realistic workload conditions, serving sensitivity, or drift aware task construction [11–13]. These are important advances in evaluation realism. The next useful step is to complement task level benchmarking with structural context divergence diagnostics that make the evaluation deployment projection more explicit, comparable, and interpretable over time.

### 7.3. Infrastructure-Aware Evaluation as a Control Surface

The serving stack is not only a delivery layer. It is a structural transformation layer that conditions how model capability is expressed in practice. Evaluation that abstracts away from serving conditions can therefore become only partially representative of production reality. In heterogeneous inference environments, execution paths vary with accelerator class, memory architecture, virtualization overhead, routing logic, and scheduling policy, creating configuration dependent performance surfaces that single context benchmarks cannot fully represent [24].

Infrastructure aware evaluation would treat these serving conditions not as external noise, but as part of the behavioural surface of the deployed system. This does not imply replacing model centric evaluation. It means extending it so that system level execution conditions become more visible as co determinants of realised behaviour. For hyperscalers, such an approach can improve comparability across deployment targets, reduce avoidable debugging effort, and support more predictable transitions from evaluation to production.

### 7.4. Cost of Structural Ignorance

Industry reporting indicates that many organisations continue to encounter unexpected cost dynamics in token based AI services as deployment complexity increases [3]. One reason is that production behaviour is shaped not only by nominal model capability, but also by retry behaviour, context expansion, orchestration overhead, and serving layer interactions that may not be visible during model selection. In such cases, evaluation deployment divergence can contribute to cost expansion even when benchmark performance appears satisfactory.

Structural visibility into evaluation deployment divergence therefore has direct operational value. It improves the likelihood that deployment decisions are made on behaviour that remains representative under production conditions rather than on results that are only locally predictive within benchmark settings. For hyperscalers, this can support earlier recognition of configuration sensitive drift, more reliable capacity planning, and a clearer understanding of how evaluation outcomes translate into realised efficiency at scale.

## 8. Conclusion

Benchmark saturation does not reduce the relevance of evaluation. It changes the level at which evaluation remains most informative. As frontier models converge across major benchmark suites, marginal score differences become less decisive as a standalone basis for production choice. Under these conditions, deployment outcomes are shaped increasingly by execution geometry, namely the interaction between model behaviour, serving stack, runtime coordination, and infrastructure context. The central implication of this paper is that evaluation and deployment can no longer be treated as adjacent but structurally equivalent stages of the same process.

The divergence analysed here is therefore not best understood as a temporary shortcoming of benchmark design. It is a structural property of large scale AI systems operating across different contexts. Evaluation environments are intentionally bounded, reproducible, and comparable. Deployment environments are variable, heterogeneous, and coupled to live runtime conditions. As serving architectures become more stateful, agentic workflows become more extended, and infrastructure paths become more differentiated, the structural distance between these two contexts becomes increasingly relevant for understanding realised system behaviour.

From this perspective, structural diagnostics do not compete with existing evaluation practice. They extend it. Their role is to make the projection from evaluation context to deployment context more explicit, more monitorable, and more interpretable over time. For hyperscalers and AI platform operators, this additional visibility can support more robust deployment selection, clearer infrastructure-model fit assessment, and more predictable translation from nominal capability to realised production performance.

The broader strategic shift is therefore straightforward. As benchmark compression increases, the question is no longer only which model performs best on a benchmark, but which model-context combination remains most coherent once deployed. In this regime, structural alignment becomes a practical differentiator, and evaluation methodology benefits from being complemented by a more explicit view of context-sensitive execution.

*When benchmarks saturate, capability ceases to be the decisive variable. Execution geometry becomes the real frontier.*

*Evaluation stability does not guarantee deployment stability.*

**Acknowledgments:** The author acknowledges prior internal architectural work that informed the conceptual development of the diagnostic perspective presented in this paper.

**Conflicts of Interest:** The author declares no conflicts of interest.

**Use of Artificial Intelligence:** Artificial intelligence tools were used for limited editorial support, including language refinement and  $\LaTeX$  formatting assistance. The scientific framing, conceptual arguments, structural taxonomy, interpretation of the literature, and all substantive research judgments were developed and verified by the author, who takes full responsibility for the content of the manuscript.

**Data Availability Statement:** No new data were generated in this study. All referenced data are available in the cited publications.

1. Maslej, N.; Fattorini, L.; Perrault, R.; Gil, Y.; Parli, V.; et al. (2025). The AI Index 2025 Annual Report. *AI Index Steering Committee, Institute for Human-Centered AI, Stanford University*. Stanford, CA.
2. Epoch AI. (2025). LLM Inference Prices Have Fallen Rapidly But Unequally Across Tasks. *Epoch AI Data Insights*. <https://epoch.ai/data-insights/llm-inference-price-trends>
3. McKinsey & Company. (2025). The State of AI in 2025: Generative AI Adoption and Cost Economics. *McKinsey Global Institute*.
4. OpenAI. (2026). Introducing GPT-5.3 Codex. *OpenAI Blog*. <https://openai.com/index/introducing-gpt-5-3-codex/>
5. Anthropic. (2026). Claude Opus 4.6. *Anthropic*. <https://www.anthropic.com/news/claude-opus-4-6>
6. Google DeepMind. (2026). Gemini 3.1 Pro Model Card. *Google DeepMind*. <https://deepmind.google/models/model-cards/gemini-3-1-pro/>
7. Alibaba Cloud. (2026). Qwen 3.5: 397B MoE with 17B Active Parameters. *Alibaba Cloud Blog*. <https://www.alibabacloud.com/blog/qwen-3-5>
8. OpenAI. (2026). OpenAI and Amazon Partnership. *OpenAI Blog*. <https://openai.com/index/amazon-partnership/>
9. Reuters. (2026). Nvidia Plans New Chip to Speed AI Processing. *Reuters Business*. <https://www.reuters.com/business/nvidia-plans-new-chip-speed-ai-processing-wsj-reports-2026-02-28/>
10. NVIDIA Developer Blog. (2026). Automating Inference Optimizations with NVIDIA TensorRT-LLM AutoDeploy. *NVIDIA Developer Blog*. <https://developer.nvidia.com/blog/automating-inference-optimizations-with-nvidia-tensorrt-llm-autodeploy/>
11. NVIDIA Research. (2026). SPEED-Bench: A Unified and Diverse Benchmark for Speculative Decoding. *NVIDIA Research Publications*. [https://research.nvidia.com/publication/2026-02\\_speed-bench-unified-and-diverse-benchmark-speculative-decoding](https://research.nvidia.com/publication/2026-02_speed-bench-unified-and-diverse-benchmark-speculative-decoding)
12. ISO-Bench Authors. (2026). ISO-Bench: Can Coding Agents Optimize Real-World Inference Workloads? *arXiv preprint*. <https://arxiv.org/html/2602.19594v1>
13. FeatureBench Authors. (2026). FeatureBench: Benchmarking Agentic Coding for Complex Feature Development. *arXiv preprint*. <https://arxiv.org/abs/2602.10975>
14. Liang, P., et al. (2023). Holistic Evaluation of Language Models. *Transactions on Machine Learning Research*.
15. Liu, X., et al. (2023). AgentBench: Evaluating LLMs as Agents. *Proceedings of ICLR 2024*.
16. Zhou, S., et al. (2023). WebArena: A Realistic Web Environment for Building Autonomous Agents. *Proceedings of ICLR 2024*. <https://webarena.dev/>

17. Goldman Sachs Research. (2025). Data Center Power Demand: The 6 Ps Driving Growth and Constraints. *GS SUSTAIN Report*.
18. International Energy Agency. (2025). Energy and AI: Energy Demand from Data Centres. *IEA Special Report*. Paris. <https://www.iea.org/reports/energy-and-ai>
19. Gartner, Inc. (2025). Predicts 2025: Agentic AI — The Evolution of Experience and Productivity. *Gartner Research Report*.
20. Galileo. (2025). The Hidden Costs of Agentic AI: Why 40% of Projects Fail Before Production. *Galileo AI Blog*. <https://galileo.ai/blog/hidden-cost-of-agentic-ai>
21. Wegener, G.H. (2026). SORT-AI: Agentic System Stability in Large-Scale AI Systems — Structural Causes of Cost, Instability, and Non-Determinism in Multi-Agent and Tool-Using Workflows. *MDPI Preprints*. doi:10.20944/preprints202601.1741.v1
22. Wegener, G.H. (2026). SORT-AI: Runtime Control Coherence in Large-Scale AI Systems. *MDPI Preprints*. doi:10.20944/preprints202601.0298.v1
23. Wegener, G.H. (2026). SORT-AI: Interconnect Stability and Cost per Performance in Large-Scale AI Infrastructure. *MDPI Preprints*. doi:10.20944/preprints202601.0161.v1
24. Wegener, G.H. (2026). SORT-AI: Accelerator Runtime Coherence in Heterogeneous AI Inference Infrastructure — A Structural Analysis of Cross-Layer Instability in Large-Scale Systems. *MDPI Preprints*. doi:10.20944/preprints202603.1614.v1