

Article

Not peer-reviewed version

SORT-AI: Structural Efficiency Recovery in Hyperscale AI Systems A Diagnostic Framework for Throughput, Control, and Orchestration Losses

[Gregor Herbert Wegener](#)*

Posted Date: 2 February 2026

doi: 10.20944/preprints202602.0015.v1

Keywords: AI infrastructure efficiency; structural diagnostics; distributed training; inference serving; agentic systems; coordination overhead; SORT framework; hyperscale systems



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a [Creative Commons CC BY 4.0 license](#), which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

SORT-AI: Structural Efficiency Recovery in Hyperscale AI Systems A Diagnostic Framework for Throughput, Control, and Orchestration Losses

Gregor Herbert Wegener 

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

Abstract

Large-scale AI infrastructure exhibits a persistent efficiency paradox: despite continuous hardware advancement and algorithmic refinement, realized performance often falls substantially below nominal capacity specifications. Contemporary optimization efforts focus primarily on accelerator utilization, model architecture, and serving algorithms, yet empirical evidence from production deployments indicates that structural losses arising from interconnect contention, control layer conflicts, and orchestration overhead constitute a distinct and largely unaddressed inefficiency surface. This work introduces the SORT-AI framework, a diagnostic approach to structural efficiency analysis in hyperscale AI systems. Rather than proposing implementation-specific optimizations, the framework formalizes the relationship between nominal hardware capacity and effective system throughput through the lens of structural coupling across three failure domains: synchronization-induced losses in distributed training, memory-control friction in inference serving, and intent propagation failures in agentic orchestration environments. We establish a taxonomy of structural losses that distinguishes between recoverable inefficiencies arising from system-level coupling and fundamental algorithmic or hardware constraints. Through application of the framework to interconnect stability, runtime control coherence, and agentic system coordination, we identify structural recovery bounds that represent efficiency improvements achievable without hardware replacement or runtime substitution. The framework operates as a diagnostic layer that surfaces inefficiencies typically obscured by component-level metrics, providing architectural visibility into coordination failures that manifest as ghost compute, stranded capacity, and emergent cost amplification. This work does not provide benchmarks, performance guarantees, or implementation prescriptions. Instead, it establishes a conceptual foundation for understanding efficiency recovery as the inversion of structural instability, offering a systems-level perspective on why optimization efforts often underperform expectations and where architectural assessment can yield operational insights.

Keywords: AI infrastructure efficiency; structural diagnostics; distributed training; inference serving; agentic systems; coordination overhead; SORT framework; hyperscale systems

1. Introduction

1.1. The Efficiency Paradox in Hyperscale AI Systems

The deployment of hyperscale AI infrastructure has accelerated rapidly over the past decade, driven by sustained advances in accelerator hardware, interconnect technologies, and large-scale training frameworks. Modern AI systems routinely operate on specialized accelerators such as NVIDIA H100-class GPUs, Google TPU generations, and domain-specific ASICs, enabling unprecedented raw computational throughput. Despite these advances, empirical analyses of production-scale AI datacenters consistently report that effective utilization remains substantially below nominal hardware capacity, often ranging between 30% and 60%.

Recent characterizations of large-scale training and inference workloads indicate that this utilization gap persists even in highly optimized environments. The discrepancy between theoretical peak performance and realized throughput suggests that performance limitations increasingly arise not from insufficient compute capability, but from system-level coordination effects that emerge as infrastructure scales in size, heterogeneity, and operational complexity.

Large foundation model training provides a particularly illustrative example. At cluster scale, coordination overheads associated with synchronization, fault handling, and topology constraints grow nonlinearly with system size. While kernel-level optimizations and architectural refinements continue to improve per-device efficiency, these improvements do not necessarily translate into proportional gains in end-to-end throughput. As a result, improvements in nominal capacity can coexist with persistent structural inefficiencies at the system level.

Complementary observations from concurrent GPU serving and pooling systems further reinforce this pattern. Even when individual accelerators operate within their expected performance envelopes, resource contention and scheduling interactions can dominate aggregate system behavior. The resulting efficiency paradox—continued hardware advancement without commensurate gains in effective capacity—indicates that coordination and coupling effects constitute a distinct analytical domain requiring architectural, rather than component-level, examination.

1.2. Utilization Gap and Coordination Dominance

The distinction between nominal capacity and effective capacity has long been recognized in warehouse-scale computing and datacenter systems. Nominal capacity reflects the theoretical maximum achievable under idealized isolation assumptions, whereas effective capacity captures realized throughput under realistic coupling conditions, including shared resources, control interactions, and operational constraints.

In contemporary AI infrastructure, coordination-dominated regimes account for a significant fraction of observed efficiency loss. In distributed training, synchronization barriers introduced by gradient aggregation and collective communication can occupy substantial portions of iteration time. During these intervals, accelerators remain powered and active, yet no forward or backward computational progress occurs. The magnitude of this effect is sensitive to network topology, interconnect bandwidth, and failure handling mechanisms, and increases with cluster scale.

Inference serving environments exhibit analogous coordination effects. Interactions between schedulers, batching strategies, and memory management—particularly under dynamic request patterns—frequently necessitate deliberate overprovisioning to maintain tail latency guarantees. Empirical reports from production systems indicate that maintaining service-level objectives often requires reserving significant excess capacity, effectively trading utilization for predictability.

At the control plane level, large-scale scheduling systems have documented that independently correct optimization decisions, when composed across multiple autonomous layers, can collectively reduce global efficiency. Such coordination losses may not manifest as explicit failures or degraded component metrics, yet they materially constrain effective throughput.

Agentic execution environments introduce an additional coordination regime. When language models operate as autonomous components with iterative planning, tool invocation, and multi-agent interaction, token consumption can increase substantially without proportional advancement toward task objectives. In these settings, utilization losses arise from orchestration dynamics rather than from core inference computation, further expanding the scope of coordination-dominated inefficiency.

1.3. Why Classical Optimization Approaches Are Insufficient

Classical performance optimization methodologies predominantly focus on localized metrics, including accelerator utilization, memory bandwidth saturation, and kernel execution efficiency. These approaches have proven highly effective for improving component-level performance and remain indispensable for low-level optimization.

However, such metrics are inherently limited in their ability to capture inefficiencies that arise from interactions among multiple autonomous components operating under partial information. Coordination losses do not necessarily manifest as degraded performance at any single layer, but instead emerge from the composition of individually optimized subsystems sharing resources and control authority.

Kernel-level optimizations improve instruction efficiency but do not address waiting time induced by synchronization or scheduling decisions external to the kernel. Similarly, model architecture refinements can reduce computational demand per inference step, yet they do not eliminate inefficiencies arising from memory contention, control plane interactions, or orchestration overhead.

Prior work on technical debt in machine learning systems has highlighted that system-level inefficiencies often remain latent until they surface as operational complexity, escalating costs, or reduced predictability. Architectural decomposition and abstraction boundaries, while essential for scalability, can inadvertently obscure coupling surfaces where efficiency losses accumulate.

These observations motivate the need for analytical tools that operate at the structural level of AI systems. Rather than replacing existing optimization techniques, such tools aim to complement them by providing visibility into coordination effects that are not observable through component-focused instrumentation alone.

1.4. Scope and Contribution

This work introduces a structural efficiency framework designed to analyze coordination-induced losses in hyperscale AI infrastructure. The scope of the framework is explicitly limited to infrastructure-level coupling effects, including synchronization barriers, control plane interactions, and orchestration dynamics arising from multi-component execution environments.

The framework does not address model-internal reasoning efficiency, prompt engineering techniques, or hardware microarchitectural design. These domains represent distinct optimization surfaces and require separate analytical and empirical treatment.

The contributions of this work are fourfold. First, we establish a formal taxonomy of structural efficiency losses, classified by their dominant coupling domain across physical synchronization, logical control, and semantic orchestration layers. Second, we formalize efficiency recovery as a process of structural inversion, distinguishing recoverable coordination losses from fundamental system constraints. Third, we map the taxonomy to diagnostic application domains spanning distributed training, inference serving, and agentic execution environments, grounding the framework in existing system classes. Fourth, we derive conservative recovery bounds indicating potential efficiency improvements achievable without hardware replacement or runtime substitution.

This framework builds on established principles in distributed systems design, multi-agent coordination, and datacenter-scale architecture. It extends recent work on agentic and tool-augmented AI systems by providing a system-level lens for understanding efficiency degradation in extended autonomous execution contexts.

2. Conceptual Foundations: The SORT-AI Framework

2.1. Nominal Capacity versus Effective Capacity

The distinction between theoretical and realized performance has been a central concern in datacenter systems design since the emergence of warehouse-scale computing [22,23]. While hardware specifications define upper performance bounds, production systems rarely operate under ideal isolation. In the context of AI infrastructure, this discrepancy becomes particularly pronounced due to tight coupling across computation, communication, and control layers.

We formalize this distinction through three capacity metrics.

Nominal Capacity (C_{nom}) denotes the theoretical maximum throughput implied by hardware specifications under idealized conditions. It assumes isolated component operation without synchro-

nization barriers, control conflicts, or orchestration delays, and thus represents an upper bound on achievable performance [14].

Effective Capacity (C_{eff}) denotes the realized throughput observed under production conditions. It incorporates synchronization overhead, control plane interactions, and orchestration effects, and is measured through end-to-end task completion metrics rather than component-local activity indicators [50].

Structural Loss (L_{struct}) captures the capacity gap arising specifically from system-level coupling effects:

$$L_{\text{struct}} = C_{\text{nom}} - C_{\text{eff}} \quad (1)$$

The corresponding utilization gap is given by:

$$U_{\text{gap}} = \frac{L_{\text{struct}}}{C_{\text{nom}}} \quad (2)$$

This formulation distinguishes structural losses from algorithmic inefficiencies or hardware constraints. Structural losses arise even when individual components operate correctly within their specifications, reflecting coordination requirements rather than fundamental computational limits [20].

2.2. Structural Efficiency and the Recovery Principle

Classical optimization approaches primarily target increases in nominal capacity through improved algorithms, enhanced hardware, or refined architectures [12,14]. Structural efficiency addresses a different objective: the reduction of L_{struct} by stabilizing coordination patterns and resolving coupling-induced conflicts.

This distinction is fundamental. Structural efficiency does not exceed nominal capacity; rather, it approaches it by eliminating preventable losses. We refer to this mechanism as the *recovery principle*. If a structural instability source S reduces effective capacity by a fraction f , then stabilization of S enables recovery of up to that same fraction f .

In distributed systems terminology, this corresponds to a shift from fault tolerance—maintaining operation despite degradation—to fault prevention, where structural sources of degradation are addressed directly [48]. The recovery principle operates through feedback mechanisms that align local decisions with global system state [47]. When scheduling, memory management, or orchestration decisions are made without sufficient cross-layer visibility, structural losses accumulate as an emergent property of distributed decision-making under incomplete information [45].

2.3. SORT-AI as a Diagnostic Framework

The Supra-Omega Resonance Theory (SORT) provides a general formalism for structural analysis across multiple domains, including AI systems and complex adaptive systems [1,2]. SORT-AI represents the application of this formalism specifically to AI infrastructure, focusing on diagnostic analysis of structural coupling effects rather than component-level optimization.

SORT-AI distinguishes three primary coupling domains.

Physical Coupling Domain. Synchronization-induced losses in distributed training arise from hardware coordination requirements, including gradient aggregation barriers, all-reduce communication patterns, and network topology constraints [7,10]. These effects are analyzed within the SORT-AI Interconnect Stability framework [3].

Logical Coupling Domain. In inference serving, memory-control friction emerges from misalignment between scheduling decisions and resource availability, particularly in KV-cache management and dynamic batching contexts [17,18]. These effects are examined through the SORT-AI Runtime Control Coherence framework [4].

Semantic Coupling Domain. In agentic orchestration, intent propagation failures arise from planning drift, redundant tool invocation, and multi-agent coordination conflicts [27,28]. These effects are addressed in the SORT-AI Agentic System Stability framework [5].

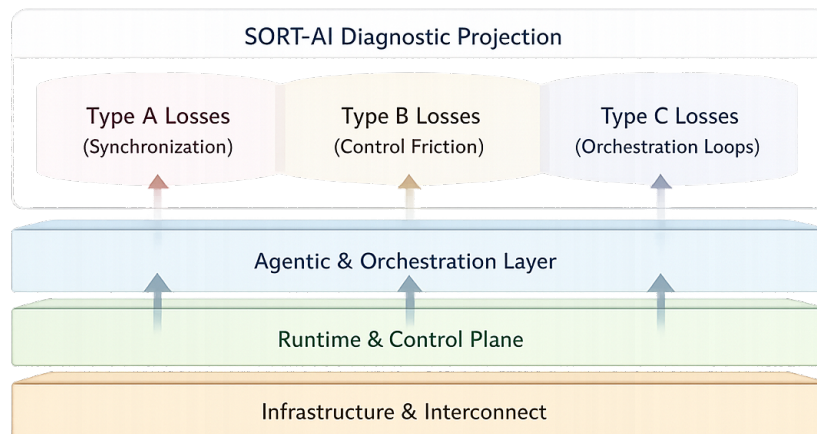


Figure 1. SORT-AI Structural Diagnostic Overview. The framework operates as a diagnostic projection across layered AI system architectures. Infrastructure, runtime control, and agentic orchestration layers are analyzed through a non-prescriptive diagnostic lens, mapping structural coupling patterns to loss categories (Type A–C) without introducing optimization or control mechanisms.

2.4. Structural Coupling and Failure Domains

A critical property of these coupling domains is their irreducibility. Addressing losses in one domain does not eliminate losses in others [46]. Optimizing interconnect topology may reduce synchronization overhead, but it does not resolve scheduling conflicts between autonomous control loops. Similarly, improvements in serving algorithms do not address divergence in agentic planning.

This irreducibility reflects fundamental properties of distributed decision-making systems. In multi-agent coordination theory, local rationality does not guarantee global optimality, and emergent system behavior can diverge substantially from component-level intent [39,40]. Complex adaptive systems exhibit emergent properties that cannot be inferred from component behavior alone [43,44].

The SORT-AI diagnostic framework therefore provides a structured methodology for decomposing system behavior into coupling-specific components, enabling targeted assessment of where structural losses accumulate and which classes of intervention are likely to yield recovery effects.

3. Taxonomy of Structural Losses

3.1. Type A: Synchronization-Induced Losses

Synchronization-induced losses arise in distributed training environments where gradient aggregation requires coordinated execution across multiple compute nodes [7,10]. During collective communication operations such as all-reduce or parameter synchronization, accelerators frequently enter idle states while awaiting completion of global coordination steps. Although hardware resources remain active and consume power, no forward or backward propagation progress occurs during these intervals, resulting in effective compute loss [8].

These losses are amplified by straggler effects in large-scale clusters. In synchronous training regimes, delays introduced by a single slow node propagate across all participating workers, forcing otherwise idle accelerators to wait [11,20]. Empirical studies of GPU datacenter workloads indicate that synchronization overhead can occupy 20–30% of iteration time in clusters exceeding several hundred accelerators [7].

Network topology further shapes the magnitude of synchronization-induced losses. Hierarchical all-reduce communication patterns generate structured traffic that can saturate interconnect links, introducing coordination delays even when individual accelerators retain computational headroom [10]. Techniques such as in-network aggregation reduce absolute synchronization time, yet the underlying coordination requirement persists as a structural constraint rather than a removable inefficiency.

The recovery vector for Type A losses therefore does not seek to eliminate synchronization, but to stabilize and localize its effects. Diagnostic analysis of gradient flow topology and interconnect

stress patterns enables identification of regimes where synchronization overhead exceeds necessary coordination costs, providing opportunities for structural recovery through improved alignment between communication patterns and system topology [3].

3.2. Type B: Memory-Control Friction Losses

Memory-control friction losses emerge predominantly in inference serving environments, where scheduling decisions are executed without complete visibility into memory resource state [17,18]. Large language model inference is characterized by attention-driven memory access patterns and persistent KV-cache growth during sequence generation. When schedulers dispatch requests based on compute availability alone, memory contention can induce latency spikes and throughput degradation [20].

Dynamic batching strategies are commonly employed to maximize throughput by grouping concurrent requests. However, batching decisions can conflict with memory management constraints as KV-cache expansion increases memory pressure. Under saturation, cached states may be evicted prematurely, forcing recomputation of previously processed tokens and reducing effective throughput. In such cases, accelerators may exhibit available compute capacity while being unable to accept new work due to memory constraints, resulting in stranded capacity.

To maintain service-level guarantees under variable load, production serving systems frequently provision 30–50% excess capacity [18,21]. When this over-provisioning arises from incomplete coordination between scheduling and memory management rather than intrinsic workload variance, it constitutes a structural loss rather than a necessary safety margin.

Large-scale scheduling systems have documented that multiple autonomous schedulers, each operating correctly within a local scope, can generate emergent conflicts that consume substantial capacity without manifesting as explicit failures [19]. Retry amplification mechanisms, where transient conflicts trigger cascaded rescheduling attempts, further contribute to capacity loss while remaining indistinguishable from normal operational behavior in aggregate metrics.

The recovery vector for Type B losses focuses on control coherence: aligning scheduling decisions with real-time memory state and resource availability [4]. This diagnostic perspective is distinct from serving algorithm optimization. Rather than improving batching heuristics, it targets the elimination of conflicts between independently rational control decisions.

3.3. Type C: Orchestration Loop Losses

Type C losses arise in agentic AI systems, where large language models function as autonomous decision-making components within extended workflows [27,28]. Unlike synchronization losses rooted in hardware coordination or control losses driven by resource contention, orchestration loop losses originate at the semantic level. They manifest as planning drift, redundant tool invocation, and failures in intent propagation across multi-agent interactions.

Characteristic patterns include abandoned planning branches that consume tokens without contributing to final outputs, repeated reasoning cycles that do not advance task completion, and tool invocations whose results are unused or contradicted by subsequent decisions [5]. These losses scale superlinearly with system complexity. As the number of agents, tools, and planning steps increases, the combinatorial expansion of possible interaction paths leads to substantial exploration that does not translate into productive outcomes [36].

Multi-agent coordination theory has long established that independently rational agents operating under partial observability can generate globally inefficient behavior [39,40]. Without stabilizing coordination mechanisms, agent interactions may interfere, producing retry cascades, duplicated work, and inconsistent state transitions [41,42].

Tool-augmented language models enable rich interaction with external systems [32,33], but this flexibility introduces additional orchestration overhead. Retrieval-augmented generation pipelines, for example, can multiply database queries and token usage through repeated embedding generation and context assembly, frequently recomputing overlapping information [6,35].

Industry analyses indicate that in recursive planning scenarios, agentic workflows can incur cost amplification exceeding two orders of magnitude relative to baseline token consumption, with effective token utilization dropping below 5% when planning loops diverge or tool results are underutilized [37, 38].

The recovery vector for Type C losses targets stabilization of intent propagation and orchestration coherence. Structural recovery operates at the coordination layer, aligning agent interactions and tool usage patterns, rather than modifying model-internal reasoning processes [5].

3.4. Secondary Loss Categories

Beyond the primary loss types, secondary structural losses arise when supporting mechanisms themselves introduce disproportionate overhead. Fault-recovery strategies can become sources of inefficiency when recovery costs exceed the impact of the faults being mitigated [48,49]. Checkpointing operations introduce synchronization barriers, retry logic amplifies coordination overhead, and roll-back procedures can cascade across dependent services [11]. In such regimes, resilience mechanisms transition from protective layers to dominant contributors to structural loss.

Retrieval-augmented generation architectures exhibit similar secondary effects. Redundant database queries, repeated embedding computation, and duplicated context assembly across parallel retrieval paths generate overlapping work that does not improve output quality [6,35].

Finally, pricing and billing models can amplify structural inefficiencies. Token-based billing magnifies the cost of abandoned reasoning paths, API-based pricing multiplies the impact of redundant tool calls, and time-based compute billing increases the economic weight of idle periods [36]. These amplification effects do not introduce new inefficiencies but increase the operational significance of existing structural losses.

4. Structural Sources Across System Classes

4.1. Distributed Training and Interconnect Effects

Large-scale foundation model training operates in a synchronization-dominated regime in which coordination overhead can exceed computational work [7,8]. The Llama 3 infrastructure report documents that gradient synchronization occupies 20–30% of iteration time in 1000+ GPU clusters, creating substantial idle periods where accelerators consume power without producing forward or backward propagation progress [8].

All-reduce operations for gradient aggregation generate hierarchical communication patterns that interact with network topology. In ring-based all-reduce implementations, bandwidth utilization approaches optimal levels, but latency accumulates linearly with cluster size [10]. Tree-based all-reduce reduces latency but introduces bottlenecks at aggregation nodes, particularly when gradient tensors exceed network switch buffer capacity.

Pipeline parallelism strategies introduce additional coordination complexity. While pipeline stages can overlap computation and communication, they create pipeline bubbles where stages idle while waiting for inputs from upstream stages [7]. The fraction of time spent in bubbles scales with pipeline depth and stage imbalance, representing structural loss independent of individual stage efficiency.

Straggler effects amplify dramatically in synchronous training contexts. A single slow accelerator in a 1024-GPU cluster can delay the entire training step, multiplying a local delay of milliseconds into aggregate idle time measured in GPU-seconds [11,20]. Mitigation strategies such as speculative execution or dynamic rescheduling reduce sensitivity to individual stragglers but introduce additional coordination overhead, trading one form of structural loss for another.

From a diagnostic perspective, recovery opportunities in distributed training lie in identifying where synchronization losses exceed necessary coordination overhead. Interconnect stress patterns, gradient flow analysis, and topology-aware scheduling can surface recovery opportunities without requiring hardware replacement or algorithm changes [3].

4.2. Inference Serving and Memory-Control Coupling

Inference serving systems operate under fundamentally different constraints than training systems. Request patterns are dynamic and unpredictable, latency requirements are strict, and resource utilization must balance throughput maximization against tail latency guarantees [17,18].

Dynamic batching aims to amortize inference overhead by processing multiple requests concurrently. Naive batching strategies that maximize batch size can violate latency service-level agreements, while conservative strategies that prioritize latency leave throughput on the table [17]. Continuous batching approaches that allow requests to join mid-sequence improve utilization but introduce scheduling complexity and potential memory contention.

Attention mechanisms in transformer architectures create characteristic memory access patterns through KV-cache management. Cache grows linearly with sequence length, and eviction policies must balance memory utilization against re-computation overhead [17]. When schedulers dispatch new requests without visibility into cache state, memory conflicts create latency spikes as cached states are evicted and later re-computed.

Production serving systems frequently provision significant excess capacity to maintain SLA compliance under dynamic load. Meta's tail utilization analysis documents that inference clusters operate at 30–50% average utilization to provide headroom for request spikes and variability [18]. While some over-provisioning is necessary, a substantial fraction arises from control plane incoherence rather than intrinsic uncertainty in request patterns [21].

Large-scale scheduling systems such as Borg demonstrate that multiple autonomous schedulers can exhibit emergent coordination failures even when individual components operate correctly [19]. Retry amplification occurs when transient conflicts trigger cascading retries that consume capacity without appearing as explicit failures.

Recovery opportunities in inference serving lie at the intersection of scheduling visibility and memory management. By surfacing memory state to schedulers and aligning batching decisions with resource availability, systems can reduce over-provisioning without compromising SLA guarantees [4].

4.3. Agentic Systems and Orchestration Overhead

The emergence of agentic AI systems represents a shift from isolated inference toward coordinated multi-step execution [27,28]. Agents decompose tasks, invoke tools, revise plans based on intermediate results, and coordinate with other agents when solving complex problems [29,30].

Tool-augmented language models enable rich interactions with external systems [32,33], but introduce orchestration overhead that scales nonlinearly with agentic complexity [34]. Agents may generate plans that are later abandoned or invoke tools whose results are unused or contradicted.

Multi-agent coordination theory shows that locally rational agents operating under partial observability do not guarantee globally efficient outcomes [39,41]. Redundant work, retry cascades, and inconsistent intermediate states emerge as structural properties of decentralized coordination.

Agentic systems exhibit characteristic ghost costs including ghost tokens, ghost planning, and ghost tool calls [5]. These losses scale superlinearly with the number of agents, tools, and planning depth [36].

Retrieval-augmented generation pipelines further amplify orchestration overhead through query multiplication, redundant embedding generation, and duplicated context assembly [6,35]. Empirical analyses indicate 3–5× query multiplication and 5–10× token overhead relative to direct generation [38].

Industry analyses suggest that orchestration costs can reach 100× baseline token consumption in recursive planning scenarios, with effective token utilization falling below 5% when planning loops diverge [36,37].

Recovery opportunities in agentic systems lie in intent coherence mechanisms that stabilize planning loops and reduce abandoned work, operating at the orchestration layer rather than within model-internal reasoning [5].

4.4. Fault-Recovery Amplification and Stability Collapse

Fault tolerance mechanisms can become primary sources of structural inefficiency when recovery overhead exceeds the cost of the faults being mitigated [48,49].

Checkpointing operations introduce synchronization barriers that pause computation [7]. When checkpoint frequency is high relative to fault probability, cumulative checkpoint overhead can exceed the expected cost of occasional restarts.

Retry mechanisms may trigger cascading retry storms when transient conflicts are not properly dampened [20]. The Borg system has documented cases where retry amplification consumed more capacity than original request processing [19].

Fault recovery can also induce cascading failures across dependent services [24]. Observability gaps further obscure recovery costs, as checkpointing and retries are aggregated into baseline performance metrics [15,16].

Diagnostic assessment must therefore distinguish between necessary resilience overhead and amplification effects. When recovery mechanisms consume more resources than the faults they mitigate, structural stability has collapsed and the resilience architecture itself requires reevaluation.

5. The Logic of Structural Inversion

5.1. Recovery as Inversion of Structural Instability

Classical optimization in computer systems primarily targets enhancement of nominal capacity through improved algorithms, more capable hardware, or refined architectures [12,14]. Structural efficiency recovery operates through a fundamentally different mechanism. Rather than increasing capability, it reduces coordination-induced losses by stabilizing structural instabilities that prevent existing capacity from being effectively realized.

The structural inversion principle formalizes this distinction. If a structural instability source S reduces effective capacity by a fraction f , then stabilization of S enables recovery of up to fraction f of that lost capacity. Recovery cannot exceed nominal capacity; it represents elimination of preventable losses rather than creation of new computational capability. In distributed systems terminology, this reframes efficiency improvement from fault tolerance—maintaining operation in the presence of failures—to fault prevention—eliminating structural conditions that generate degradation [48].

The inversion mechanism operates through feedback processes that align local decision-making with global system state [47]. When schedulers operate without visibility into resource saturation, when agents plan without awareness of parallel execution conflicts, or when synchronization barriers form without topology awareness, structural losses emerge as an aggregate property of distributed decision-making under incomplete information [45,46].

Structural stabilization does not eliminate coordination requirements. Distributed training continues to require synchronization, inference serving continues to require scheduling, and agentic systems continue to require orchestration. However, by surfacing coupling patterns and aligning decisions with actual system state, systems can approach coordination-theoretic lower bounds rather than operating with excess overhead [39].

5.2. Why Recovery Bounds Are Conservative

Recovery bounds are intentionally conservative because not all observed inefficiency arises from structural coupling. Performance gaps may also result from algorithmic limitations, hardware constraints, or irreducible computational complexity [50,51]. The framework explicitly targets only the structural component of observed losses.

For example, consider a distributed training system operating at 50% effective utilization. Detailed analysis may attribute 20% of the gap to synchronization overhead (structural, Type A), 15% to data pipeline limitations (algorithmic), 10% to kernel inefficiencies, and 5% to memory bandwidth saturation. Structural recovery can address only the synchronization component; the remaining losses require orthogonal interventions.

Even when structural losses are correctly identified, recovery mechanisms introduce their own overhead. Visibility instrumentation consumes resources, stabilization logic adds control complexity, and improvements in one coupling domain may introduce secondary coupling surfaces elsewhere [25]. Conservative bounds reflect these implementation realities by discounting theoretical recovery potential to account for practical limitations.

Recovery potential also varies substantially across architectures and workloads. A training cluster with balanced topology may exhibit lower synchronization loss than aggregate statistics suggest, while an inference system with predictable demand may require less over-provisioning than systems exposed to bursty traffic [23]. Accordingly, the bounds presented here represent indicative ranges rather than guaranteed outcomes.

5.3. Throughput Recovery versus Cost Reduction

Structural recovery manifests in two operationally distinct but complementary modes: throughput recovery and cost reduction. Both represent efficiency improvements, but they address different system objectives and arise in different coupling regimes [52].

Throughput Recovery Mode. In synchronization-dominated regimes such as distributed training, stabilization reduces idle time and re-execution overhead, increasing effective throughput on fixed hardware. Training clusters that spend 20–30% of iteration time in synchronization barriers can recover 5–15% of that idle fraction through interconnect stabilization, topology-aware scheduling, and gradient flow optimization [3,7]. This yields more useful computation per GPU-hour without additional hardware investment.

Cost Reduction Mode. In orchestration-dominated regimes such as agentic systems, stabilization eliminates non-productive token consumption, redundant tool invocations, and abandoned planning branches. Multi-agent workflows that exhibit extreme token amplification can recover 10–25% of ghost token costs through intent coherence mechanisms and orchestration-level alignment [5,36]. This enables equivalent task completion with substantially reduced resource expenditure.

Certain loss categories exhibit dual recovery characteristics. In inference serving, control coherence simultaneously reduces over-provisioning and increases request throughput. By aligning scheduling decisions with memory availability, systems can serve more requests per accelerator while maintaining service-level guarantees [4,18].

Table 1. Indicative Structural Recovery Bounds by Application Domain

Application	Loss Type	Primary Mode	Conservative Bound
ai.01 Interconnect Stability	Type A	Throughput	5–15% effective throughput
ai.04 Control Coherence	Type B	Cost + Throughput	5–15% ghost cost elimination
ai.13 Agentic Stability	Type C	Cost	10–25% token cost reduction
ai.17 Fault-Recovery Stability	Amplification	Ghost Compute	5–10% recovery overhead

These bounds are derived from structural loss analysis rather than controlled benchmarks. They characterize recoverable loss envelopes under favorable diagnostic conditions.

5.4. Validity Conditions

Structural recovery operates under explicit validity conditions that delineate the scope of applicability.

No Hardware Changes Required. Recovery mechanisms operate at the coordination layer and do not require accelerator upgrades, memory expansion, or interconnect replacement. Structural efficiency reduces L_{struct} rather than increasing C_{nom} .

No Runtime Substitution Required. The framework provides diagnostic insight rather than prescriptive replacement. Existing schedulers, orchestrators, and serving stacks remain in place, reducing operational disruption and avoiding defensive reactions from production teams [25].

System Scale and Complexity Dependence. Structural losses amplify with scale and coupling complexity. Large training clusters and multi-agent workflows exhibit greater recovery potential than small or loosely coupled systems [28,31].

Workload Sensitivity. Synchronization-intensive workloads benefit from Type A recovery, control-heavy serving workloads from Type B recovery, and coordination-dominated agentic workloads from Type C recovery. Systems dominated by algorithmic or hardware constraints may exhibit limited structural recovery opportunity.

Observability Requirements. Structural diagnostics require visibility into idle time, retry behavior, token consumption, and coordination overhead. Aggregate throughput and latency metrics alone are insufficient [51].

No Performance Guarantees. Recovery bounds are indicative rather than contractual. The framework provides a methodology for assessment, not a guarantee of outcome. System-specific analysis is required to validate recovery potential [50].

6. Diagnostic Application Mapping: SORT-AI Core

6.1. Core-3 Primary Diagnostic Layer

The SORT-AI framework maps to three primary diagnostic applications, each addressing a distinct coupling domain and loss taxonomy category. These applications provide structured entry points for architectural assessment rather than prescriptive implementation guidelines.

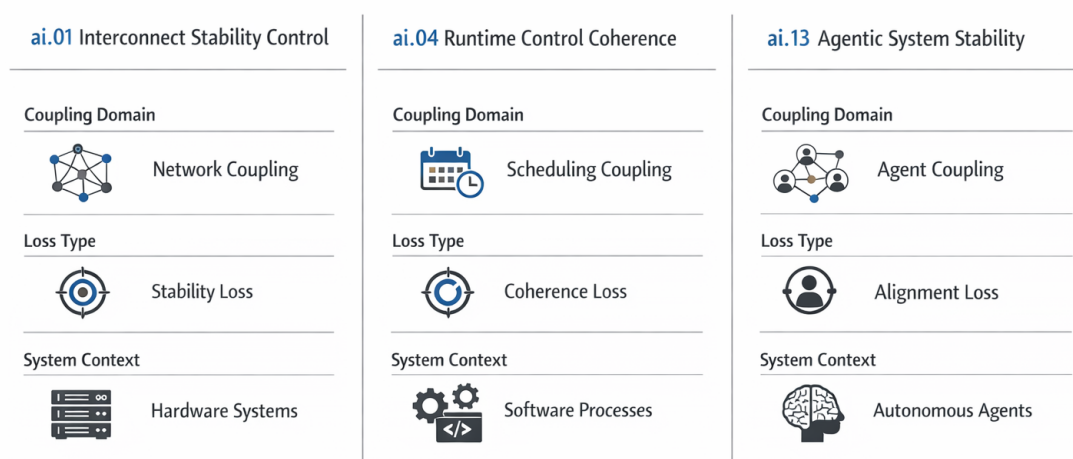


Figure 2. Core-3 Diagnostic Application Mapping. The three primary SORT-AI diagnostic applications each correspond to a dominant coupling domain and structural loss category. The mapping illustrates analytical scope rather than prescriptive intervention, providing structured entry points for architectural assessment.

Application ai.01: Interconnect Stability Control

Target: Type A losses (synchronization-induced) in distributed training contexts.

Domain: Large-scale foundation model training, model parallelism, and gradient aggregation systems [7,8].

Diagnostic Focus: Analysis of synchronization patterns, all-reduce communication bottlenecks, straggler effects, and network topology interactions. The diagnostic framework surfaces conditions under which coordination overhead exceeds necessary minimums and identifies topology-specific vulnerability patterns [3,10].

Recovery Vector: Throughput recovery via structural stabilization of gradient flow patterns and reduction of synchronization-induced idle time. Conservative bounds indicate 5–15% effective throughput improvement without hardware modification.

Applicability: Training clusters exceeding approximately 64 GPUs with measurable synchronization overhead. Recovery potential amplifies with cluster size and network depth [7].

Application ai.04: Runtime Control Coherence

Target: Type B losses (memory-control friction) in inference serving systems.

Domain: LLM inference serving, dynamic batching architectures, and scheduler–memory interactions [17,18].

Diagnostic Focus: Analysis of scheduling decision alignment with resource availability, KV-cache contention patterns, retry amplification, and control-plane incoherence. The diagnostic framework identifies emergent conflicts generated by multiple autonomous schedulers and exposes resource-visibility gaps that induce unnecessary over-provisioning [4,19].

Recovery Vector: Dual recovery mechanism combining cost reduction through ghost-cost elimination (5–15%) and throughput recovery through reduced over-provisioning, enabling higher utilization while preserving SLA compliance [21].

Applicability: Serving systems with dynamic load patterns, multiple scheduling layers, and explicit SLA requirements. Effects amplify in multi-tenant environments and under high request variance [20].

Application ai.13: Agentic System Stability

Target: Type C losses (orchestration loops) in agentic AI systems.

Domain: Multi-agent frameworks, tool-calling workflows, RAG systems with agentic control, and autonomous planning environments [27,28].

Diagnostic Focus: Analysis of planning drift, tool invocation efficiency, ghost-token accumulation, and intent propagation across agent interactions. The diagnostic framework surfaces cases where planning loops diverge without termination criteria and where tool outputs are systematically underutilized [5, 31].

Recovery Vector: Cost reduction through ghost-token elimination (10–25%), stabilization of planning loops, and optimization of tool invocation. Particularly effective in recursive planning scenarios where orchestration complexity amplifies overhead [36,37].

Applicability: Agentic systems with multiple agents (>3), rich tool ecosystems (>5 tools), or recursive planning patterns. Effects amplify superlinearly with agentic complexity [34].

6.2. Secondary Diagnostic Extensions

While Core-3 applications address approximately 70–75% of identified structural inefficiency themes, targeted extensions expand diagnostic coverage to roughly 90% through focused assessment of specific coupling patterns.

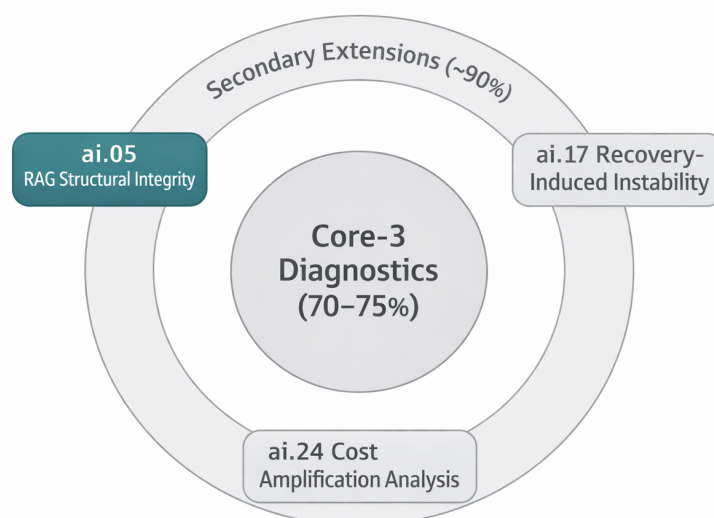


Figure 3. Secondary Diagnostic Extensions Coverage. Core-3 applications capture the majority of recurring structural inefficiency patterns. Targeted extensions expand diagnostic coverage by addressing specific coupling mechanisms that emerge in retrieval, recovery, and cost-amplification contexts.

Extension ai.05: RAG Structural Integrity

Addresses: Database query multiplication, embedding redundancy, and context assembly inefficiency in retrieval-augmented generation systems [6,35].

Recovery Focus: 5–10% orchestration efficiency improvement through reduction of redundant database queries (3–5× multiplication) and embedding overhead (5–10× token inflation) [38].

Priority: High for RAG-heavy deployments, moderate otherwise.

Extension ai.17: Recovery-Induced Instability

Addresses: Checkpointing overhead, retry amplification, and cascading failures triggered by fault-tolerance mechanisms [11,48].

Recovery Focus: 5–10% ghost-compute elimination when recovery mechanisms dominate operational overhead. The diagnostic identifies when fault-tolerance architecture has transitioned from a resilience layer to a primary inefficiency source.

Priority: High for systems with frequent checkpointing or aggressive retry policies.

Extension ai.24: Cost Amplification Analysis

Addresses: Emergent cost patterns arising from interactions between pricing structures and structural inefficiencies. Token-based billing magnifies ghost-token costs, while API invocation pricing multiplies tool-calling overhead [36].

Recovery Focus: Cost-containment visibility rather than direct recovery. The diagnostic surfaces locations where pricing models amplify structural losses, enabling targeted intervention [52].

Priority: High for cost-sensitive deployments, particularly agentic systems with elevated token consumption.

6.3. Application Clusters and Market Positioning**AI Diagnostic Applications**

Figure 4. Application Clusters and Organizational Alignment. Diagnostic applications are grouped into functional clusters aligned with operational responsibilities and architectural domains. The clustering reflects organizational interfaces rather than analytical separation within the SORT-AI framework.

Diagnostic applications are organized into functional clusters aligned with operational team structures and architectural domains.

Cluster A: Infrastructure & Interconnect

Target Teams: Training infrastructure engineers, HPC system architects, distributed systems specialists.

Primary Applications: ai.01 (Interconnect Stability Control).

Extensions: ai.06 (Energy–Interconnect Coupling), ai.08 (Scalability Certification).

Characteristic Challenge: Synchronization overhead in large-scale training environments [7,8].

Cluster C: Runtime & Control

Target Teams: Inference serving teams, orchestration platform developers, scheduling system engineers.

Primary Applications: ai.04 (Runtime Control Coherence).

Extensions: ai.05 (RAG Integrity), ai.09 (Control-Flow Stability), ai.17 (Recovery Stability).

Characteristic Challenge: Control-plane conflicts in inference-serving systems [17,19].

Cluster D: Agentic & Emergent

Target Teams: Multi-agent framework developers, agentic system operators, LLM application architects.

Primary Applications: ai.13 (Agentic System Stability).

Extensions: ai.24 (Cost Amplification).

Characteristic Challenge: Orchestration overhead in autonomous workflows [27,31].

Applications function as diagnostic entry points rather than product specifications. They provide structured methodologies for identifying where structural losses accumulate and which recovery mechanisms are likely to yield efficiency improvements. This positioning avoids prescriptive implementation mandates while enabling targeted architectural assessment [25].

7. Scope, Non-Claims, and Validity Conditions

7.1. What This Framework Does Not Address

The framework explicitly excludes model-level efficiency concerns including internal LLM reasoning optimization, chain-of-thought prompting strategies, and model architecture design [53,54]. These represent distinct optimization surfaces operating at the algorithmic layer rather than the infrastructure coordination layer. Improving reasoning efficiency within a model does not address synchronization losses in distributed training or control conflicts in serving systems.

The framework does not address algorithmic training improvements such as gradient compression techniques, mixed-precision strategies, or optimizer selection [14]. These approaches target enhancement of C_{nom} through better algorithms. Structural efficiency targets reduction of L_{struct} through coordination improvement. Both are necessary; neither is sufficient alone.

The framework operates independently of hardware microarchitecture details including GPU design, memory hierarchy specifications, or interconnect topology as an implementation prescription [23]. It provides diagnostic visibility into coupling patterns without mandating specific hardware configurations. Recovery mechanisms are architecture-agnostic, operating at the coordination layer above hardware specifics.

The framework maintains vendor neutrality, avoiding prescriptions for CUDA kernel tuning, TPU compilation strategies, or cloud provider-specific configurations. This positioning ensures applicability across heterogeneous infrastructure environments and avoids creating vendor lock-in through framework adoption.

7.2. What This Framework Does Not Provide

The framework provides diagnostic visibility and structural analysis but does not offer benchmarks, performance comparisons, or vendor rankings [50]. It surfaces inefficiency patterns and indicates recovery potential without prescribing specific performance targets or guaranteeing measurement outcomes.

Conservative recovery bounds (5–25% depending on loss type) represent indicative ranges derived from structural analysis, not contractual guarantees [51]. Individual systems may exhibit recovery potential outside these ranges depending on workload characteristics, architectural details, and implementation quality. The framework provides assessment methodology rather than turnkey results.

The framework operates as a diagnostic layer that surfaces structural losses without specifying fixes [25]. It provides architectural visibility into coupling failures but leaves implementation decisions to engineering teams familiar with specific system constraints and operational requirements. This positioning respects existing expertise while enabling informed decision-making.

The framework does not replace existing runtimes, schedulers, or orchestration systems. It operates alongside current infrastructure as a visibility layer that surfaces inefficiencies typically obscured by component-level metrics [46]. This complementary positioning avoids triggering defensive responses from teams operating production systems.

7.3. Validity Conditions for Framework Application

The framework applies most effectively to hyperscale deployments where coordination overhead becomes visible. Training clusters exceeding 64 GPUs, inference systems with hundreds of concurrent requests, and agentic workflows with multiple agents and tool ecosystems exhibit structural losses at levels where diagnostic assessment yields actionable insights [7,23]. Small-scale systems (single-node training, low-traffic serving, single-agent workflows) may exhibit negligible structural losses, limiting recovery opportunity.

The framework targets coordination-dominated workloads where coupling overhead represents a substantial fraction of total execution time. Synchronization-intensive workloads (distributed training with frequent all-reduce operations), control-heavy workloads (inference serving under dynamic load with SLA requirements), and orchestration-dominated workloads (agentic systems with recursive planning) demonstrate characteristic structural loss patterns [7,17,27]. Workloads limited by algorithmic complexity or hardware constraints rather than coordination will see limited benefit.

Structural diagnostics require instrumentation beyond standard throughput and latency metrics. Visibility into idle time distributions, retry rate patterns, token consumption breakdowns, and coupling-specific overhead is necessary for decomposing performance gaps into structural versus non-structural components [15,51]. This represents an investment in observability infrastructure as a prerequisite for framework application.

General recovery bounds do not substitute for system-specific analysis. Individual architectures require targeted assessment to validate recovery potential, identify dominant loss patterns, and determine applicability of stabilization mechanisms [50]. The framework provides methodology and classification structure rather than universal quantitative predictions.

7.4. Relationship to Other Optimization Approaches

Structural efficiency is complementary to, not competitive with, existing optimization approaches [12,13]. Kernel optimization improves C_{nom} and should continue. Model architecture refinement enhances inference quality and should continue. Structural efficiency addresses L_{struct} and provides additive benefit when component-level optimization has been exhausted.

Traditional optimization targets throughput, latency, and FLOPS as primary metrics [50]. Structural efficiency targets idle time, ghost costs, retry overhead, and coordination losses—metrics that often remain invisible in aggregate performance dashboards [15]. Both metric spaces are necessary for comprehensive system optimization; neither is sufficient alone.

The framework suggests a staged optimization approach: First, optimize kernels and models to establish best-case C_{nom} [14]. Second, profile system behavior to identify L_{struct} through coupling-specific instrumentation [51]. Third, apply structural recovery mechanisms to reduce L_{struct} toward coordination-theoretic minimums. This staged approach ensures that optimization efforts target appropriate layers sequentially rather than conflating distinct efficiency opportunities.

8. Conclusions

8.1. Framework Summary

This work has formalized the efficiency paradox in hyperscale AI infrastructure: despite continuous hardware advancement and algorithmic refinement, realized performance often falls substantially below nominal capacity specifications [7,23]. We have established that structural losses arising from synchronization overhead, control conflicts, and orchestration inefficiencies constitute a distinct failure mode requiring analysis at the system architecture level rather than the component level.

The SORT-AI framework provides three primary contributions. First, a formal taxonomy of structural losses classified by coupling domain: Type A (synchronization-induced), Type B (memory-control friction), and Type C (orchestration loops) [3–5]. Second, formalization of the efficiency recovery principle as structural inversion—reducing L_{struct} by stabilizing coordination patterns rather than enhancing C_{nom} through algorithmic or hardware improvements [47]. Third, mapping of the taxonomy to diagnostic applications across training, inference, and agentic system classes, with conservative recovery bounds (5–25%) indicating efficiency improvement potential without hardware replacement.

A central insight is that structural losses remain largely invisible to component-level metrics. Standard instrumentation capturing accelerator utilization, memory bandwidth, and throughput cannot distinguish between algorithmic constraints and coupling-induced inefficiencies [50,51]. The framework operates as a diagnostic layer that decomposes performance gaps into structural versus non-structural components, enabling targeted intervention where coordination overhead dominates system behavior.

8.2. Implications for Hyperscale AI Architecture

For infrastructure teams operating large-scale training clusters, the framework suggests that synchronization overhead represents a first-order architectural concern rather than a secondary optimization target [7,10]. Coordination-aware scheduling, topology-sensitive gradient flow, and interconnect stress analysis provide complementary visibility to hardware monitoring, surfacing recovery opportunities that remain hidden in aggregate utilization metrics.

For system architects designing inference serving platforms, the framework highlights that control coherence is a structural requirement for efficient operation at scale [17,19]. When multiple autonomous schedulers operate without visibility into shared resource state, emergent conflicts consume capacity without surfacing as explicit failures. Providing schedulers with memory state visibility and aligning batching decisions with resource availability represents a distinct optimization surface from serving algorithm refinement.

For AI researchers developing agentic systems, the framework indicates that intent propagation failures and orchestration overhead require semantic-layer diagnostics beyond model capability assessment [27,28]. Ghost token accumulation, planning loop divergence, and tool invocation redundancy represent structural inefficiencies that persist even when underlying models demonstrate strong reasoning capabilities [5]. Efficiency in agentic systems requires coordination mechanisms that stabilize planning and align agent interactions, complementing ongoing model improvements.

8.3. Future Directions

The framework establishes conceptual foundations and diagnostic methodology but does not provide empirical validation through production measurements. Controlled experiments isolating structural components from algorithmic and hardware constraints represent a natural next step [55]. Production deployments offer validation opportunities where structural analysis predictions can be compared against observed recovery effects under controlled intervention.

The framework currently addresses distributed training, inference serving, and agentic orchestration. Extensions to additional domains including multi-modal systems (vision-language coordination), federated learning (coordination across trust boundaries), and hybrid cloud deployments (cross-datacenter coupling) represent opportunities for broadening diagnostic coverage [53].

Integration with existing frameworks represents a practical deployment path. Observability platforms (Prometheus, Datadog), orchestration systems (Kubernetes, Ray), and cost management tools (FinOps platforms) could incorporate structural visibility as a complementary diagnostic layer [13]. This integration would enable operational teams to leverage structural diagnostics without adopting entirely new tooling.

Theoretical work connecting structural coupling analysis to distributed systems theory, control theory, and multi-agent coordination offers opportunities for formal validation of recovery bounds and applicability conditions [39,45,47]. Coordination-theoretic minimums for various system classes would provide principled targets for recovery assessment rather than empirically-derived conservative bounds.

8.4. Closing Remarks

This framework does not solve efficiency challenges. It provides a structured approach to understanding why efficiency targets often underperform expectations and where architectural assessment can yield operational insights. Structural efficiency recovery represents a diagnostic opportunity complementary to ongoing optimization efforts in hardware, algorithms, and architectures.

Systems that achieve nominal capacity at the component level yet underperform at the system level may benefit from structural visibility that surfaces coordination losses typically obscured by conventional metrics [15,46]. The efficiency paradox persists not because individual components fail to perform well, but because their interactions create emergent inefficiencies that remain invisible until formalized and measured.

The SORT-AI framework invites continued exploration of structural coupling as a first-order architectural concern in hyperscale AI systems. As infrastructure scales and system complexity increases, coordination overhead will occupy an expanding fraction of total operational cost [23,36]. Understanding, measuring, and addressing these structural losses represents a necessary evolution in how we approach efficiency in the era of hyperscale AI deployment.

Acknowledgments: The author acknowledges the open sharing of infrastructure insights and operational experiences by engineering teams at Meta, Google, Databricks, NVIDIA, and other organizations whose documentation of production challenges informed this analysis.

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

Use of Artificial Intelligence: The author confirms that no artificial intelligence tools were used in the generation of scientific concepts, theoretical frameworks, or results. Automated tools were used exclusively for editorial language refinement and L^AT_EX formatting assistance.

References

1. Wegener, G. H. (2025). SORT-AI: A Projection-Based Structural Framework for AI Safety—Alignment Stability, Drift Detection, and Scalable Oversight. *Preprints* **2024121334**. DOI:10.20944/preprints202412.1334.v2
2. Wegener, G. H. (2025). SORT-CX: A Projection-Based Structural Framework for Complex Systems—Operator Geometry, Non-Local Kernels, Drift Diagnostics, and Emergent Stability. *Preprints* **2024121431**. DOI:10.20944/preprints202412.1431.v1
3. Wegener, G. H. (2025). SORT-AI: Interconnect Stability and Cost per Performance in Large-Scale AI Infrastructure—A Structural Analysis of Runtime Instability in Distributed Systems. *Preprints* **2026010161**. DOI:10.20944/preprints202601.0161.v1
4. Wegener, G. H. (2025). SORT-AI: Runtime Control Coherence in Large-Scale AI Systems—Structural Causes of Cost, Instability, and Non-Determinism Beyond Interconnect Failures. *Preprints* **2026010298**. DOI:10.20944/preprints202601.0298.v1
5. Wegener, G. H. (2025). 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. *Preprints* **2026011741**. doi:10.20944/preprints202601.1741.v1

6. Wegener, G.H. (2025). SORT-AI: A Structural Safety and Reliability Framework for Advanced AI Systems with Retrieval-Augmented Generation as a Diagnostic Testbed. *Preprints* **2024121345**. DOI:10.20944/preprints202412.1345.v1
7. Jeon, M., Venkataraman, S., Phanishayee, A., et al. (2024). Characterization and Prediction of Deep Learning Workloads in Large-Scale GPU Datacenters. *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC '24)*.
8. Dubey, A., Jauhri, A., Pandey, A., et al. (2024). The Llama 3 Herd of Models. *arXiv preprint arXiv:2407.21783*.
9. Zhai, E., et al. (2025). Aegaeon: Effective GPU Pooling for Concurrent LLM Serving on the Market. *Proceedings of the 29th ACM Symposium on Operating Systems Principles (SOSP '25)*.
10. Rabinovich, E., et al. (2023). Scaling Distributed Machine Learning with In-Network Aggregation. *Proceedings of NSDI'23*, 1–19.
11. Ananthanarayanan, G., et al. (2013). Effective Straggler Mitigation: Attack of the Clones. *Proceedings of NSDI'13*, 185–198.
12. Delimitrou, C., & Kozyrakis, C. (2014). Quasar: Resource-Efficient and QoS-Aware Cluster Management. *Proceedings of ASPLOS'14*, 127–144. DOI:10.1145/2541940.2541941
13. Zaharia, M., et al. (2018). Accelerating the Machine Learning Lifecycle with MLflow. *IEEE Data Engineering Bulletin* **41**(4), 39–45.
14. Patterson, D. A., & Hennessy, J. L. (2016). *Computer Organization and Design: The Hardware/Software Interface* (5th ed.). Morgan Kaufmann. ISBN 978-0-12-407726-3.
15. Sculley, D., et al. (2015). Hidden Technical Debt in Machine Learning Systems. *Proceedings of NeurIPS 2015*, 2503–2511.
16. Amershi, S., et al. (2019). Software Engineering for Machine Learning: A Case Study. *Proceedings of ICSE-SEIP'19*, 291–300. DOI:10.1109/ICSE-SEIP.2019.00042
17. Databricks Engineering. (2024). LLM Inference Performance Engineering: Best Practices. *Databricks Engineering Blog*. databricks.com/blog
18. Meta Engineering. (2024). Taming Tail Utilization of Ads Inference at Meta Scale. *Meta Engineering Blog*. engineering.fb.com
19. Tirmazi, M., et al. (2020). Borg: The Next Generation. *Proceedings of the 15th European Conference on Computer Systems (EuroSys '20)*, 1–14. DOI:10.1145/3342195.3387517
20. Dean, J., & Barroso, L. A. (2013). The Tail at Scale. *Communications of the ACM* **56**(2), 74–80. DOI:10.1145/2408776.2408794
21. Kannan, R. S., et al. (2019). GrandSLAM: Guaranteeing SLAs for Jobs in Microservices Execution Frameworks. *Proceedings of EuroSys'19*, 1–16. DOI:10.1145/3302424.3303958
22. Barroso, L. A., Clidaras, J., & Hölzle, U. (2013). The Datacenter as a Computer: An Introduction to the Design of Warehouse-Scale Machines (2nd ed.). *Synthesis Lectures on Computer Architecture* **8**(3), 1–154. DOI:10.2200/S00516ED2V01Y201306CAC024
23. Barroso, L. A., Hölzle, U., & Ranganathan, P. (2019). The Datacenter as a Computer: Designing Warehouse-Scale Machines (3rd ed.). *Synthesis Lectures on Computer Architecture* **13**(3), 1–189. DOI:10.2200/S00874ED3V01Y201809CAC046
24. Gunawi, H. S., et al. (2014). What Bugs Live in the Cloud? A Study of 3000+ Issues in Cloud Systems. *Proceedings of SoCC'14*, 1–14. DOI:10.1145/2670979.2670986
25. Ousterhout, J. (2018). *A Philosophy of Software Design*. Yaknyam Press. ISBN 978-1-7321022-0-0.
26. Uber Engineering. (2024). Vertical CPU Scaling: Reduce Cost of Capacity and Increase Reliability. *Uber Engineering Blog*. uber.com/blog
27. Xi, Z., et al. (2023). The Rise and Potential of Large Language Model Based Agents: A Survey. *arXiv preprint arXiv:2309.07864*.
28. Wang, L., et al. (2024). A Survey on Large Language Model Based Autonomous Agents. *Frontiers of Computer Science* **18**(6), 186345. DOI:10.1007/s11704-024-40231-1
29. Yao, S., et al. (2023). ReAct: Synergizing Reasoning and Acting in Language Models. *Proceedings of ICLR 2023*.
30. Shinn, N., et al. (2023). Reflexion: Language Agents with Verbal Reinforcement Learning. *Proceedings of NeurIPS 2023*.
31. Wu, Q., et al. (2023). AutoGen: Enabling Next-Gen LLM Applications via Multi-Agent Conversation. *arXiv preprint arXiv:2308.08155*.

32. Schick, T., et al. (2023). Toolformer: Language Models Can Teach Themselves to Use Tools. *Proceedings of NeurIPS 2023*.
33. Patil, S. G., et al. (2023). Gorilla: Large Language Model Connected with Massive APIs. *arXiv preprint arXiv:2305.15334*.
34. Qin, Y., et al. (2023). Tool Learning with Foundation Models. *arXiv preprint arXiv:2304.08354*.
35. Lewis, P., et al. (2020). Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. *Proceedings of NeurIPS 2020*.
36. IDC & DataRobot. (2025). The Hidden AI Tax: Cost Control in the Age of GenAI and Agentic Workflows. *IDC Market Spotlight*.
37. Stevens Institute of Technology. (2025). The Hidden Economics of AI Agents: Managing Token Costs and Latency Trade-offs. *Stevens Online Blog*. online.stevens.edu/blog
38. Nous Research. (2025). Token Efficiency Across Language Models. *Technical Report*.
39. Wooldridge, M. (2009). *An Introduction to MultiAgent Systems* (2nd ed.). John Wiley & Sons. ISBN 978-0-470-51946-2.
40. Dorri, A., Kanhere, S. S., & Jurdak, R. (2018). Multi-Agent Systems: A Survey. *IEEE Access* **6**, 28573–28593. DOI:10.1109/ACCESS.2018.2831228
41. Panait, L., & Luke, S. (2005). Cooperative Multi-Agent Learning: The State of the Art. *Autonomous Agents and Multi-Agent Systems* **11**(3), 387–434. DOI:10.1007/s10458-005-2631-2
42. Stone, P., & Veloso, M. (2000). Multiagent Systems: A Survey from a Machine Learning Perspective. *Autonomous Robots* **8**(3), 345–383. DOI:10.1023/A:1008942012299
43. Holland, J. H. (1995). *Hidden Order: How Adaptation Builds Complexity*. Addison-Wesley. ISBN 978-0-201-44230-4.
44. Levin, S. A. (2003). Complex Adaptive Systems: Exploring the Known, the Unknown and the Unknowable. *Bulletin of the American Mathematical Society* **40**(1), 3–19. DOI:10.1090/S0273-0979-02-00965-5
45. Lamport, L. (1978). Time, Clocks, and the Ordering of Events in a Distributed System. *Communications of the ACM* **21**(7), 558–565. DOI:10.1145/359545.359563
46. Saltzer, J. H., Reed, D. P., & Clark, D. D. (1984). End-to-End Arguments in System Design. *ACM Transactions on Computer Systems* **2**(4), 277–288. DOI:10.1145/357401.357402
47. Hellerstein, J. L., Diao, Y., Parekh, S., & Tilbury, D. M. (2004). *Feedback Control of Computing Systems*. John Wiley & Sons. ISBN 978-0-471-26637-2.
48. Avizienis, A., et al. (2004). Basic Concepts and Taxonomy of Dependable and Secure Computing. *IEEE Transactions on Dependable and Secure Computing* **1**(1), 11–33. DOI:10.1109/TDSC.2004.2
49. Perrow, C. (1984). *Normal Accidents: Living with High-Risk Technologies*. Basic Books. ISBN 978-0-465-05142-9.
50. Lilja, D. J. (2000). *Measuring Computer Performance: A Practitioner's Guide*. Cambridge University Press. ISBN 978-0-521-64105-4.
51. Jain, R. (1991). *The Art of Computer Systems Performance Analysis*. John Wiley & Sons. ISBN 978-0-471-50336-1.
52. Brynjolfsson, E., & Hitt, L. M. (2000). Beyond Computation: Information Technology, Organizational Transformation and Business Performance. *Journal of Economic Perspectives* **14**(4), 23–48. DOI:10.1257/jep.14.4.23
53. Bommasani, R., et al. (2021). On the Opportunities and Risks of Foundation Models. *arXiv preprint arXiv:2108.07258*.
54. Wei, J., et al. (2022). Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. *Proceedings of NeurIPS 2022*.
55. Kapoor, S., & Narayanan, A. (2024). Leakage and the Reproducibility Crisis in ML-Based Science. *Patterns* **5**(4), 100804. DOI:10.1016/j.patter.2023.100804
56. Google Research. (2024). Solving Virtual Machine Puzzles: How AI is Optimizing Cloud Computing. *Google Research Blog*.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.