

# Resolving the Starvation Paradox: Efficient Disaggregated LLM Serving with Phasor

Anonymous ACL submission

## Abstract

Phase-disaggregated LLM serving separates compute-bound prefill from memory-bound decode to improve efficiency and isolation. We show that physical separation alone does not guarantee either: under stochastic production bursts, schedulers driven by phase-local proxies overload compute-saturated prefill workers, throttle the ready rate into decode, and leave KV-resident decode GPUs idle—an efficiency failure rather than a tail-latency artifact.

We propose Virtual Freeness ( $\Phi$ ), a phase-consistent, dimensionless headroom signal that makes prefill compute saturation and decode effective KV capacity loss directly comparable, enabling principled cross-phase trade-offs while remaining orthogonal to existing execution substrates. Built on  $\Phi$ , Phasor is a holistic control framework for disaggregated clusters that performs micro-scale routing, meso-scale decode maintenance, and macro-scale pool re-balancing.

Evaluated on representative NLP tasks and realistic bursty workloads, Phasor reduces end-to-end P99 latency by 25.7–38.9% compared to state-of-the-art baselines. Crucially, by suppressing fragmentation from a peak of 49.7% in baselines to a 6.2% average, it preserves effective KV capacity, thereby boosting SLO-compliant concurrency per-GPU by 1.1–2.0 $\times$ , translating to a 12–50% reduction in per-request serving costs. We will make our code and data publicly available upon acceptance

## 1 Introduction

Autoregressive LLM serving exhibits a fundamental resource duality: prefill is compute-bound (Zhao et al., 2024; Heo et al., 2024; Agrawal et al., 2024; Jha et al., 2024), whereas decode is memory-bound (Hong et al., 2024a; Kamath et al., 2025; Shazeer, 2019). A growing line of work (Zhou et al., 2024; Zhuang et al., 2024; Hong et al., 2024b; Zhong et al., 2024; Ding and Yang, 2025; Patel

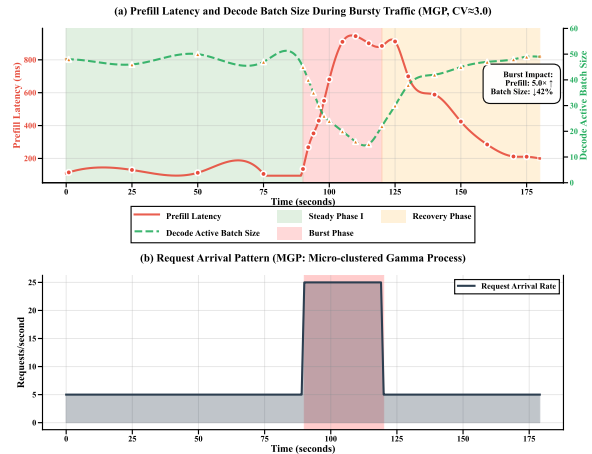


Figure 1: **Pipeline starvation under disaggregation.** Prefill-side congestion throttles the decode-ready rate, driving a high-latency, low-concurrency regime where KV-resident decode GPUs idle.

et al., 2024; Qin et al., 2024) advocates *phase disaggregation* by physically disaggregating prefill and decode to improve serving efficiency and pursue performance isolation.

However, we show that physical separation alone guarantees neither isolation nor efficiency. As Figure 1 shows, stochastic bursts trigger upstream stalls that starve the decode phase, collapsing active batch sizes. This traps the system in a *high-latency, low-concurrency* regime where decode GPUs sit idle despite holding expensive resident KV states—an efficiency failure, not merely a latency artifact.

We attribute this to a *control-plane disconnect*: existing systems rely on phase-local, incomparable proxies (e.g., queue length (Patel et al., 2024)). This blinds schedulers to compute saturation, causing them to overload workers that appear available under memory-centric views, thereby amplifying head-of-line blocking and starving downstream resources.

To bridge this, we propose Virtual Freeness ( $\Phi$ ),

a unified phase-consistent scheduling interface projecting heterogeneous constraints into a dimensionless headroom signal. By rendering prefill compute saturation and decode effective capacity loss directly comparable,  $\Phi$  enables principled cross-phase tradeoffs. Crucially, this interface is *orthogonal* to execution substrates (e.g., planned phase splits (Zhong et al., 2024), streaming (Feng et al., 2025), or disaggregated KV memory (Qin et al., 2024)) equipping their control plane with a common language for cross-phase decisions.

Built on  $\Phi$ , we present Phasor, a holistic control framework for disaggregated clusters that integrates dimensionless load modeling with predictive resource orchestration. Phasor collects lightweight telemetry, projects each worker state into  $\Phi$ , and applies actions at multiple time scales: micro-scale request routing, meso-scale decode-side maintenance, and macro-scale pool rebalancing.

Across ShareGPT-derived NLP tasks (Chiang et al., 2023) and Azure bursts workload (Wang et al., 2025b), Phasor consistently improves end-to-end tail latency and *SLO-constrained efficiency*: at 40–60 req/s, it cuts P99 latency by 38.9% and 25.7% compared to baselines. Subject to strict SLOs, it boosts compliant concurrency per-GPU by 1.1–2.0 $\times$  (reducing GPU cost per request by 12–50%). It further suppresses fragmentation (49.65% peak in baselines  $\rightarrow$  6.20% average), preserving effective KV capacity and sustaining concurrency.

Our contributions are:

- We identify a control-plane failure mode in disaggregated serving where misaligned phase-local metrics induce *pipeline starvation*, yielding high latency and low effective concurrency despite phase separation.
- We propose Virtual Freeness ( $\Phi$ ), a unified scheduling interface that provides a phase-consistent headroom signal to directly trade off prefill compute saturation against decode-side effective capacity loss.
- We build Phasor, a holistic control framework for disaggregated clusters driven by  $\Phi$ , and demonstrate improved end-to-end performance and SLO-constrained efficiency on diverse NLP and realistic bursty workloads.

## 2 Motivation

The transition from monolithic to disaggregated architectures is predicated on a core promise: by

physically decoupling the compute-bound prefill phase from the memory-bound decode phase, systems can eliminate resource contention and achieve performance isolation (Zhong et al., 2024; Hu et al., 2024b; Jiang et al., 2025).

We show this promise can fail even under strict physical separation, when the control plane relies on misaligned scheduling signals: upstream congestion can still propagate downstream and collapse end-to-end performance and serving efficiency.

**Controlled setup.** We mirror DistServe’s physical phase separation but intentionally remove system-specific heuristics, creating a clean setting. This neutral setup (8x NVIDIA V100, NVLink within each phase pod, Qwen2.5-7B-Instruct FP16) ensured that any observed bottlenecks are attributable to the control signal rather than architectural noise.

### 2.1 The Reality: The Starvation Paradox

We inject a BurstGPT-derived shock–recovery trace (Wang et al., 2025b) using an MGP arrival model ( $CV \approx 3.0$ ) to emulate micro-clustered bursts ( $\lambda=5 \rightarrow 25, \text{req/s}$  at  $t=90\text{--}120, \text{s}$ ). Each request features an average input length of 1024 tokens and an output length of 256 tokens, mirroring RAG-like workloads.

As shown in Figure 1, a burst triggers a clear *scissor effect*: prefill latency rises sharply (from  $\sim 116$  ms to  $>900$  ms), while the decode-side active batch size drops from  $\sim 48$  to  $<16$ .

This is *pipeline starvation*: overloaded prefill workers stall upstream, draining the decode-ready queue and collapsing the active decode batch, which lowers effective utilization and SLO-compliant concurrency per GPU. Decode GPUs still keep KV states resident, but idle behind upstream head-of-line blocking while waiting for prefill outputs. This is an efficiency failure rather than a mere tail-latency artifact, showing that physical phase separation alone does not guarantee performance isolation

### 2.2 The Diagnosis: The Utilization Phantom

Why does the scheduler allow prefill to become the blocking bottleneck? The root cause stems from a structural blind spot: single-dimensional metrics fail to capture the conflicting resource demands of compute-bound prefill versus memory-bound decode.

We verify this by profiling an isolated prefill pool

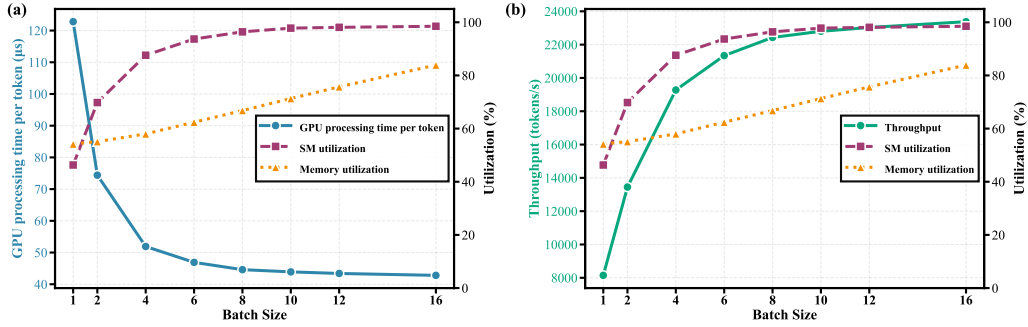


Figure 2: **The Utilization Phantom.** While memory utilization (blue) suggests the prefill node is underloaded ( $\approx 60\%$ ), compute units (SM, orange) are already fully saturated ( $> 93\%$ ). Schedulers relying on memory metrics are blind to this saturation.

(Figure 2). As batch size scales, computational efficiency improves (per token processing time drops  $122.8\mu\text{s} \rightarrow 42.8\mu\text{s}$ ) but quickly hits a hard wall: SM utilization saturates ( $> 93.7\%$  at batch size 6).

Crucially, however, memory utilization hovers at  $\sim 60\%$ . Figure 2(b) confirms that throughput gains diminish sharply beyond batch size 8 (only 21% improvement), despite 38% of memory capacity sitting idle.

This discrepancy creates a *Utilization Phantom*: the worker appears *available* in memory space while being *saturated* in compute time. Schedulers relying solely on memory headroom are thus blind to this compute saturation, erroneously routing burst traffic to overloaded nodes and triggering the downstream starvation observed in Section 2.1.

### 2.3 Consequence: From Skew to Stall

The utilization phantom manifests as queue divergence in the prefill pool: a memory-centric policy repeatedly routes bursts to memory-available yet compute-saturated workers, creating stragglers. In our burst trace, the busiest prefill worker builds a 40+ queue while others stay around  $\sim 15$ , whereas a compute-aware policy keeps queues near-uniform (details in Appendix A). This skew amplifies upstream head-of-line blocking, reduces the decode-ready rate, and pushes the system into a high-latency, low-utilization regime, hurting both tail performance and cost efficiency despite phase disaggregation.

### 2.4 The Control Plane Disconnect

The observed queue divergence exposes a systemic disconnect: while the *mechanisms* of disaggregation have matured, the *metrics* governing them remain dangerously fragmented.

Prior research has rigorously established the

physical foundations for efficiency: Phase Disaggregation (Zhong et al., 2024) decouples interference, State Disaggregation (Qin et al., 2024) pools global capacity, and Stream Scheduling (Feng et al., 2025) introduces pipeline fluidity. These works brilliantly solve *where* to place computations and *how* to move data.

### The Missing Link: Metric Incomparability.

However, the *control signals* driving these decisions operate in silos. Consequently, each system’s control plane optimizes a different local proxy metric, e.g., *token velocity* (Feng et al., 2025) or *queue length* (Patel et al., 2024). Crucially, they lack a common denominator to weigh the quadratic compute cost of a prefill burst against the linear memory risk of a decode stream.

### Our Solution: The Phasor Framework.

To bridge this semantic gap, we identify the need for a unified scheduling interface. We first propose Virtual Freeness ( $\Phi$ ), a dimensionless projection designed to translate distinct physical constraints (Time vs. Space) into a unified language. Crucially, this interface is orthogonal to existing execution paradigms, equipping the control plane with a bilingual capability that transforms chaotic interference into a tractable optimization objective for cross-phase control.

Built upon this interface, we propose Phasor, a holistic control framework that integrates dimensionless load modeling with predictive resource orchestration.

## 3 Methodology

We now formalize the unified interface, and show how Phasor uses it to drive closed-loop routing and maintenance in disaggregated clusters.

### 3.1 Virtual Freeness: The Unified Mathematical Form

To optimize across heterogeneous phases, we formalize the system state using a unified mathematical framework. For any worker  $w$ , Virtual Freeness ( $\Phi$ ) is a *SLO-normalized residual headroom*:

$$\Phi(w) = 1 - \frac{\mathcal{U}_{virt}(w)}{\mathcal{C}_{phys}(w)}. \quad (1)$$

Here,  $\mathcal{C}_{phys}$  is the hard physical budget and  $\mathcal{U}_{virt}$  is the predicted virtual load under the corresponding phase constraint.

Crucially, while the *form* of  $\Phi$  is unified, its *internal definition* adapts dynamically to the phase-specific physics of the worker. We formulate this duality using the following system of equations:

$$\Phi(w) = \begin{cases} w \in \mathcal{P}(\text{Time-Bound}): \\ 1 - \frac{\sum_{r \in \mathcal{Q}_w} (\alpha l_{in}(r)^2 + \beta \cdot P_m \cdot l_{in}(r))}{\mathbb{P}_{peak} T_{SLO}^{TTFT}} \\ w \in \mathcal{D}(\text{Space-Bound}): \\ 1 - \frac{N_{used}(w) + \mathcal{H}_{frag}(w)}{\mathbb{M}_{total}} \end{cases} \quad (2)$$

For the Prefill Phase,  $\mathcal{U}$  represents the Virtual Computational Load, defined as the estimated compute cost for all requests in the current instance queue  $\mathcal{Q}$ , where  $l_{in}$  is the input length of request  $r$ . The quadratic term  $\alpha \cdot l_{in}^2$  captures the dominant self-attention scaling, while the linear term  $\beta \cdot P_m \cdot l_{in}$  captures FFN-related compute, where  $P_m$  is the (active) parameter count obtained from model metadata.  $\mathcal{C}_{\mathcal{T}}$  denotes the compute budget over the TTFT window, determined by the hardware peak throughput at the deployed precision ( $\mathbb{P}_{peak}$ ) and the TTFT budget  $T_{SLO}^{TTFT}$ . We discuss the generalization and calibration of this time-bound surrogate in Appendix B.3.

For the Decode Phase, the capacity  $\mathcal{C}$  is defined by the hard memory limit  $\mathbb{M}_{total}$ . Here, the load includes the count of physically allocated KV blocks  $N_{used}$  and a fragmentation penalty  $\mathcal{H}_{frag}$ , defined as stranded capacity on worker  $w$ —free blocks that are physically available but practically unusable for large bursts without consolidation.

All hyperparameters and system constants (e.g.,  $\alpha, \beta, \mathbb{P}_{peak}, T_{SLO}^{TTFT}, \mathbb{M}_{total}, \mathcal{H}_{frag}$ ), along with cal-

### Algorithm 1 Phasor Holistic Control Loop

---

**Require:** Request stream  $R$ , pools  $\mathcal{P}, \mathcal{D}, \mathcal{M}$   
**Require:** Calibrated params  $\alpha, \beta, \theta_{frag}, \theta_{flow}$

- 1: **while** system running **do**
- Micro: Gradient Routing** ( $t \sim \mu s$ )
- 2: **if** new request  $r$  arrives **then**
- 3:    $w^* \leftarrow \arg \max_{w \in \mathcal{P}} \Phi_{\mathcal{T}}(w)$
- 4:   **Dispatch**  $r$  to  $w^*$
- 5: **end if**
- Meso: Proactive De-fragmentation** ( $t \sim ms$ )
- 6: **for all** decode worker  $d \in \mathcal{D}$  **do**
- 7:   compute  $\Phi_S(d) \triangleright$  space-bound freeness, Eq. (2)
- 8:   **if**  $\Phi_S(d) < \theta_{frag}$  **and**  $\mathcal{H}_{frag}(d)$  is high **then**
- 9:     **Trigger** ASYNCDERAG( $d$ )
- 10:   **end if**
- 11: **end for**
- Macro: Elastic Flow Control** ( $t \sim s$ )
- 12:    $\bar{\Phi}_{\mathcal{P}} \leftarrow \text{AVGFREENESS}(\mathcal{P});$
- 13:    $\bar{\Phi}_{\mathcal{D}} \leftarrow \text{AVGFREENESS}(\mathcal{D})$
- 14:   **if**  $\bar{\Phi}_{\mathcal{P}} < \theta_{flow} \cdot \bar{\Phi}_{\mathcal{D}}$  **and**  $\mathcal{M} \neq \emptyset$  **then**
- 15:      $m \leftarrow \mathcal{M}.\text{POP}();$  **Promote**  $m$  to  $\mathcal{P}$
- 16:   **else if**  $\bar{\Phi}_{\mathcal{D}} < \theta_{flow} \cdot \bar{\Phi}_{\mathcal{P}}$  **and**  $\mathcal{M} \neq \emptyset$  **then**
- 17:      $m \leftarrow \mathcal{M}.\text{POP}();$  **Promote**  $m$  to  $\mathcal{D}$
- 18:   **end if**
- 19:   **Sleep** for control interval  $\Delta t$
- 20: **end while**

---

ibration and some parameters' sensitivity analyses, are provided in Appendix B and F.

### 3.2 Phasor Architecture

Phasor is a centralized control plane over three pools: Prefill ( $\mathcal{P}$ ), Decode ( $\mathcal{D}$ ), and an elastic Mixed pool ( $\mathcal{M}$ ) (Figure 3). It collects lightweight telemetry (per-worker queue metadata and decode-side KV allocation/fragmentation maps), projects each worker state into the unified  $\Phi$  space, and enforces decisions at multiple time scales: micro request routing, meso decode maintenance, and macro pool rebalancing.

### 3.3 Actuation: Multi-Scale Control Mechanisms

Equipped with the unified  $\Phi$  signal, Phasor operationalizes the control loop through three synergistic mechanisms spanning different time scales. Visualization of scheduling mechanisms is illustrated in Appendix C.

$\Delta t$  denotes the sampling interval. We conducted a sensitivity analysis on  $\Delta t$  in the Appendix F.

**Zero-Downtime KV Migration.** Phasor pipelines KV migration by overlapping asynchronous KV replication with ongoing inference and using a brief final handoff for consistency; detailed mechanisms and overhead in Appendix D.

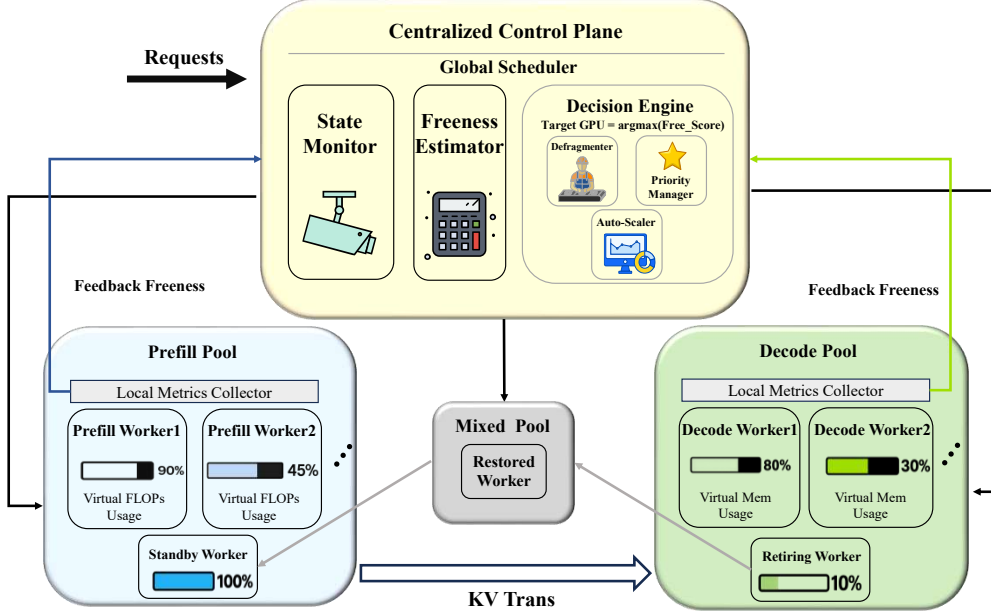


Figure 3: **Overview of the Phasor Architecture.** The Global Scheduler orchestrates resources via a closed loop: collecting raw telemetry, projecting it into the unified  $\Phi$  space, and executing actuation decisions.

**Micro (Routing).** Phasor dispatches each arrival to the prefill worker with the highest time-bound freeness  $\Phi_{\mathcal{T}}$ . This greedy move-to-the-best- $\Phi$  rule steers bursty, compute-heavy requests away from workers already dominated by quadratic prefill costs, mitigating the queue divergence described in Section 2.

**Meso (De-fragmentation).** Phasor treats fragmentation as stranded capacity and triggers ASYNCDEFRAG when  $\Phi_S(d) < \theta_{frag}$ . Rather than performing physical compaction (obviated by PagedAttention), it asynchronously migrates a subset of KV states to other workers. This effectively performs online bin-packing, consolidating scattered capacity on the source node to accommodate new heavy requests. (details in Appendix D).

**Macro (Flow Control).** Phasor compares the mean freeness of  $\mathcal{P}$  and  $\mathcal{D}$  and promotes a standby node from  $\mathcal{M}$  when one pool’s mean freeness drops below the other by a factor of the *flow-control ratio*  $\theta_{flow}$ . The choice and sensitivity of  $\theta_{flow}$  are detailed in Appendix B.5 and F.

## 4 Experiments

Our evaluation aims to answer three fundamental questions:

**RQ1 (End-to-End Performance):** Does Phasor suppress cross-phase interference and deliver consistently lower end-to-end tail latency across NLP-critical workloads and real-world production

bursts?

**RQ2 (Efficiency):** Can Phasor substantially amplify the serving efficiency of fixed hardware budgets under strict SLO constraints?

**RQ3 (Validity):** Are gains attributable to Phasor’s key mechanisms, and are results stable to threshold/sampling variations?

### 4.1 Experimental Setup

**Testbed and Models.** We conduct all experiments on a server equipped with  $8 \times$  NVIDIA V100-32GB GPUs interconnected via NVLink, representative of a standard high-performance inference pod. This setting leaves limited KV-cache headroom, making capacity sensitive to fragmentation. We deploy Qwen2.5-7B-Instruct (FP16) as the primary model. We also employ Llama-3-8B for ablation studies to verify cross-model generalization.

**Workloads: Semantic Task Archetypes.** To bridge the gap between system metrics and NLP semantics, we do not rely on random synthetic traces. Instead, we construct three distinct workload scenarios by sampling directly from the ShareGPT dataset (Chiang et al., 2023): First, Context-Heavy Tasks feature long inputs and short outputs. This pattern dominates prefill compute time, stressing the  $O(L^2)$  attention bottleneck. Second, Generation-Heavy Tasks utilize short inputs but long outputs. This pattern stresses memory capacity and is highly susceptible to KV fragmenta-

ID	In ( $L_{in}$ )	Out ( $L_{out}$ )	NLP Task
Sharegpt-L-S	~2048	~128	Summarization
Sharegpt-S-L	~128	~2048	Code Generation
BurstGPT	1024 (avg)	256 (avg)	Realistic Mix
Sharegpt-S-S*	128	128	Short QA
Sharegpt-L-L*	2048	512	Translation
Sharegpt-M-M*	1024	1024	Dialog

Table 1: Evaluated Workloads.

tion. Finally, For production bursts, we use BurstGPT (Wang et al., 2025b) derived from Azure traces to reflect real-world arrival variability. Table 1 summarizes these characteristics.

**SLO settings and metrics.** We report TTFT, TPOT, and end-to-end latency, and evaluate serving efficiency via SLO-constrained capacity under a fixed 8-GPU budget. SLO thresholds (TTFT/TPOT) are configured per workload scenario; the exact settings are summarized in the Appendix B.7.

**Baselines.** We compare Phasor against two state-of-the-art systems under an identical 8-GPU budget. DistServe (Zhong et al., 2024) represents the static disaggregation paradigm. Splitwise (Patel et al., 2024) represents standard disaggregation with a reactive Mixed Pool, routing requests based on queue length. These baselines cover the dominant static vs. reactive control paradigms in phase-disaggregated serving.

## 4.2 End-to-End Performance (RQ1)

To answer RQ1, we evaluate Phasor’s latency performance under heterogeneous bottlenecks. Figure 4 presents the latency panorama across three critical NLP task archetypes. Comprehensive results for all workloads are provided in Appendix E.

### Scenario 1: Realistic Production (BurstGPT).

Across all burst phases, Phasor reduces the average total requests P99 latency by 34.8% compared to Splitwise and 22.4% compared to DistServe. Under high load (40–60 req/s), reactive signals mask the phase-specific bottleneck, causing cascading stragglers. Virtual Freeness exposes this risk at dispatch, routing requests to workers with residual headroom to flatten load variance. Phasor therefore bounds total requests P99 latency at 206.8 s, versus 338.6 s (Splitwise) and 278.4 s (DistServe), a 38.9% and 25.7% reduction.

### Scenario 2: Summarization – Prefill Saturation.

Across all arrival rates, Phasor reduces average

TTFT P99 by 44.2% versus Splitwise and 34.1% versus DistServe. Phasor’s gain comes from making the *superlinear* long-context prefill cost a first-class routing signal: by pricing an  $l_{in}^2$ -scaled cost in  $\Phi_{\mathcal{T}}$ , it avoids straggler-dominated queues where a few long requests dominate prefill service time and inflate TTFT, preventing head-of-line blocking from cascading into decode starvation.

### Scenario 3: Generation – Memory Fragmentation.

Across all rates, Phasor reduces average TPOT by 15.6%/7.6% vs. Splitwise/DistServe. Specifically at  $\lambda=40$ , it bounds TPOT P99 at 112.8,ms (vs. 136.8/121.2,ms) and total requests P99 latency at 174.6,s (vs. 216.4/192.4,s). These gains stem from preemptive memory hygiene guided by Virtual Freeness. While baselines lose effective capacity to fragmentation, causing false-OOMs that throttle concurrency, Phasor proactively defragment memory. This preserves contiguous space, ensuring stable decode concurrency and low TPOT.

## 4.3 Resource Efficiency & Mechanism (RQ2)

We quantify serving efficiency by *SLO-constrained capacity*: the largest closed-loop *in-flight concurrency* that a fixed 8-GPU cluster can sustain while keeping both TTFT P99 and TPOT P99 within their workload-specific SLO budgets. We report the exact ( $T_{SLO}^{TTFT}, T_{SLO}^{TPOT}$ ) settings per workload in Table 3 (Appendix B.7).

### Capacity Gains

Figure 5 demonstrates that Phasor achieves superior SLO-compliant concurrency compared to baselines across diverse NLP workloads on an 8-GPU budget. Specifically, it delivers  $1.85\times$  capacity over DistServe ( $1.99\times$  vs. Splitwise) on the realistic *Production Mix*. These gains persist across distinctive regimes:  $1.66\times$  on memory-bound *Code Generation*,  $1.41\times$  on dual-bottleneck *Translation*, and  $1.13$ – $1.89\times$  across *Dialog*, *Summarization*, and *Q&A*. Overall, Phasor improves SLO-compliant concurrency per GPU by  $1.1$ – $2.0\times$ . Under a standard closed-loop assumption where the completion rate scales approximately with sustainable concurrency at fixed client think-time and workload mix, this corresponds to cutting the GPU cost per SLO-satisfying request by roughly 12–50% on the same budget.

To understand why Phasor unlocks more SLO-compliant concurrency, we examine the decode-side memory state and quantify fragmentation over time.

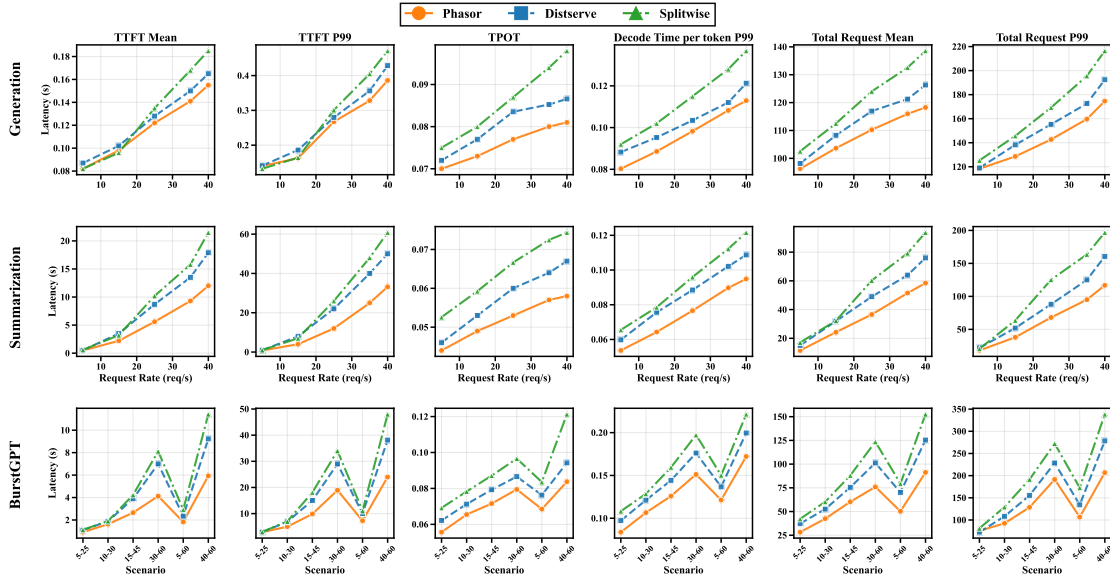


Figure 4: **End-to-end Latency Performance.** Phasor establishes a superior Pareto frontier compared to baseline.

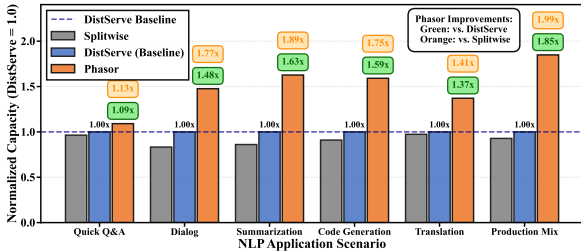


Figure 5: **SLO-constrained capacity.** Phasor yields the largest gains in generation-heavy and bursty workloads under the same 8-GPU budget.

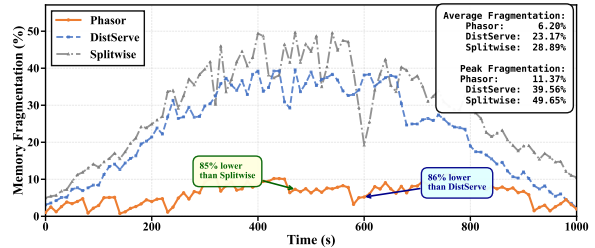


Figure 6: **Mechanism Verification: Fragmentation Suppression.** Baselines accumulate fragmentation and peak at 49.65%, while Phasor bounds fragmentation at a 6.20% average

447 **Attribution: Proactive De-fragmentation.** Fig-  
 448 ure 6 shows that baselines accumulate severe frag-  
 449 mentation (49.7%) under bursts. While PagedAt-  
 450 tention handles physical non-contiguity, it cannot  
 451 consolidate free blocks scattered across workers.  
 452 Guided by Eq. (2), Phasor proactively triggers  
 453 ASYNCDEFRAG to asynchronously migrate frag-  
 454 mented memory between workers, effectively re-  
 455 claiming contiguous capacity. Consequently, Pha-  
 456 sor maintains a low 6.2% average fragmentation  
 457 under the same load.

458 By reducing peak fragmentation at  $t=470$  s by  
 459 85.3% (7.29% vs 49.65%), Phasor converts wasted  
 460 holes into usable slots, confirming proactive hy-  
 461 giene as the decisive lever for the capacity gains in  
 462 Figure 5.

463 **Energy Efficiency.** Finally, we verify that these  
 464 capacity gains translate to sustainability. By amor-  
 465 tizing static power over higher throughput, Phasor  
 466 reduces the energy cost per token by 28.0% com-

pared to the Splitwise baseline (0.175 J/T vs. 0.243  
 J/T); see Appendix G for the detailed analysis.

#### 4.4 Architectural Validity (RQ3)

470 To attribute Phasor’s gains to its design (rather  
 471 than incidental heuristics), we conduct subtrac-  
 472 tive ablations by disabling (i) Virtual Freeness  
 473 ( $\Phi$ ), (ii) proactive de-fragmentation, and (iii) the  
 474 Mixed Pool. Table 2 reports P99 latency and *SLO-*  
 475 *constrained capacity* normalized to DistServe.

476 **Virtual Freeness is the primary driver of tail**  
 477 **latency.** Removing  $\Phi$  causes the largest and  
 478 most consistent P99 regression across workloads.  
 479 On Qwen2.5-7B-Instruct, Summarization TTFT  
 480 P99 degrades from 0.66 to 0.92 (+39%), and  
 481 Production-Mix total P99 degrades from 0.74 to  
 482 0.90 (+22%), showing that the unified time-bound  
 483 signal is essential to avoid straggler-dominated pre-  
 484 fill queues and the resulting HOL blocking.

Table 2: **Ablation of Phasor components.** Normalized to DistServe (Baseline=1.00). Removing  $\Phi$ , Mixed Pool, or De-frag degrades performance significantly.

System Variant	Norm. Latency P99 ( $\downarrow$ )			Norm. Capacity ( $\uparrow$ )		
	Summ. (TTFT)	Code (Total)	Prod. Mix (Total)	Code Gen.	Summ.	Prod. Mix
<i>Data Format: Qwen2.5-7B-Instruct (Left) / Llama-3-8B (Right)</i>						
Splitwise (Static)	1.22 / 1.24	1.13 / 1.15	1.34 / 1.36	0.91 / 0.90	0.86 / 0.85	0.93 / 0.92
DistServe (Baseline)	1.00 / 1.00	1.00 / 1.00	1.00 / 1.00	1.00 / 1.00	1.00 / 1.00	1.00 / 1.00
<b>Phasor (Full System)</b>	<b>0.66 / 0.69</b>	<b>0.91 / 0.93</b>	<b>0.74 / 0.76</b>	<b>1.59 / 1.52</b>	<b>1.63 / 1.58</b>	<b>1.85 / 1.78</b>
<i>Ablations (Impact of removing components)</i>						
– w/o Virtual Freeness	0.92 / 0.94	0.95 / 0.97	0.90 / 0.92	1.42 / 1.36	1.55 / 1.50	1.52 / 1.46
– w/o Mixed Pool	0.72 / 0.75	0.93 / 0.95	0.82 / 0.84	1.50 / 1.44	1.58 / 1.53	1.48 / 1.42
– w/o De-frag	0.68 / 0.71	0.96 / 0.98	0.82 / 0.84	1.35 / 1.30	1.60 / 1.55	1.68 / 1.62

**De-fragmentation is decisive for memory-bound capacity.** Disabling de-fragmentation yields the largest capacity drop in generation-heavy workloads: Code-Generation capacity falls from 1.59 to 1.35 (-15%). This suggests that  $\Phi$ -based routing is necessary to avoid placing load on risky workers, but de-fragmentation is needed to reclaim fragmentation-induced *effective* KV-cache capacity and keep decode concurrency stable.

**The Mixed Pool specializes in burst elasticity.** Removing the Mixed Pool primarily hurts burst capacity: Production-Mix capacity drops from 1.85 to 1.48 (-20%), while its impact on steady synthetic regimes is smaller, confirming that the Mixed Pool mainly provides macro-scale elasticity rather than per-request routing benefits.

**Robustness to parameter drift.** We sweep  $\theta_{frag}$ ,  $\theta_{flow}$ , and the control epoch  $\Delta t$  in Appendix F and observe stable performance over a broad range, indicating that Phasor does not rely on fragile tuning.

**Cross-model generalization.** The same *specialization pattern* holds on Llama-3-8B (Table 2):  $\Phi$  remains the dominant driver of tail-latency improvements across workloads, while de-fragmentation is the key enabler for capacity in memory-bound generation, and the Mixed Pool primarily affects Production-Mix capacity.

## 5 Related Work

**Phase Disaggregation.** Splitwise (Patel et al., 2024) and DistServe (Zhong et al., 2024) formalized separating prefill/decode phases, employing specialized pools with queue-based routing (queue-/token-backlog proxies) or offline profiling for planning-based split under TTFT/TPOT constraints. P/D-Serve (Jin et al., 2024) extends this with dynamic mapping for large-scale deploy-

ments.

**Dynamic & Stream-Based Scheduling.** Recent work enhances execution fluidity. WindServe (Feng et al., 2025) introduces stream-based scheduling to minimize bubbles. TaiChi (Wang et al., 2025a) dynamically toggles aggregation levels based on SLO regimes, while DéjàVu (Strati et al., 2024) leverages KV-cache streaming as a substrate for prompt-token disaggregation and fast recovery under failures.

**State-Centric Optimization.** Addressing the memory wall, Mooncake (Qin et al., 2024) and MemServe (Hu et al., 2024a) propose disaggregated architectures pooling global resources (GPU HBM and CPU DRAM). semi-PD (Hong et al., 2025) advocates storage aggregation, while Nexus (Shi et al., 2025) explores fine-grained intra-engine isolation.

## 6 Conclusion

This work demonstrates that physical separation alone fails to guarantee performance isolation. We identify a control-plane gap: phase-local, misaligned metrics can overload compute-saturated prefill workers, throttling the ready-task rate into decode so decode GPUs sit idle while still keeping KV states resident. To bridge this, we propose Phasor, a holistic control framework built upon the Virtual Freeness interface. By translating conflicting Time (latency) and Space (fragmentation) constraints into a unified coordinate system, Phasor transforms the control plane from a set of disjoint heuristics into a tractable gradient optimization problem, effectively unlocking the true efficiency potential of disaggregated serving.

## 555 Limitations

556 The limitations below primarily affect the *actuation*  
557 *cost* and *scalability* of the framework, rather than  
558 the phase-consistent scheduling interface itself.

559 **State-transfer bandwidth and topology.** Phasor’s  
560 elastic rebalancing relies on amortizing KV/state  
561 transfers (Patel et al., 2024; Zhong et al., 2024;  
562 Qin et al., 2024); when the interconnect is low-  
563 bandwidth (e.g., Ethernet-only or PCIe-constrained  
564 fabrics), migration overhead can contend with serv-  
565 ing traffic and dilute the net capacity gains (Patel  
566 et al., 2024; Zhong et al., 2024; Qin et al., 2024). In  
567 such regimes, Phasor naturally degrades to a purely  
568 routing-based controller (using  $\Phi$  for dispatch), as  
569 the cost of physical rebalancing would outweigh  
570 the fragmentation benefits.

571 **Control-plane scalability and fault tolerance.**  
572 We prioritize global scheduling optimality, which  
573 fits pod-scale deployments typical of high-  
574 performance inference. While sufficient for  
575 our evaluation scope, a single global sched-  
576 uler may become a bottleneck at data-center  
577 scales (Schwarzkopf et al., 2013; Ousterhout et al.,  
578 2013). Future work involves extending Phasor to a  
579 hierarchical architecture, where local agents handle  
580 micro-second actuation while a global coordinator  
581 manages coarse-grained pool sizing.

582 **Extremely long contexts.** For workloads with  
583 contexts orders of magnitude larger than our test  
584 set (e.g., >100K tokens), the sheer cost of moving  
585 KV states dominates, regardless of bandwidth (Pa-  
586 tel et al., 2024; Kwon et al., 2023). In these sce-  
587 narios, Phasor’s interface should be paired with  
588 copy-avoidant mechanisms (e.g., remote memory  
589 disaggregation or logical caching) (Gu et al., 2017;  
590 Huang et al., 2022; Kwon et al., 2023; Qin et al.,  
591 2024) to achieve rebalancing via *access redirec-*  
592 *tion* rather than physical state movement (Gu et al.,  
593 2017; Huang et al., 2022).

## 594 References

595 Amey Agrawal and 1 others. 2024. Taming (throughput-  
596 latency) tradeoff in llm inference with sarathi-serve.  
597 *18th USENIX OSDI*.

598 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng,  
599 Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan  
600 Zhuang, Yonghao Zhuang, Joseph E. Gonzalez,  
601 Ion Stoica, and Eric P. Xing. 2023. Vicuna:  
602 An open-source chatbot impressing gpt-4 with  
603 90%\* chatgpt quality. [https://lmsys.org/blog/  
604 2023-03-30-vicuna/](https://lmsys.org/blog/2023-03-30-vicuna/). Accessed: 2025-02-20.

Zhiqiang Ding and Tongkai Yang. 2025. Dynamicat-  
tention: Dynamic kv cache for disaggregate llm in-  
ference. In *ICASSP 2025-2025 IEEE International  
Conference on Acoustics, Speech and Signal Process-  
ing (ICASSP)*, pages 1–5. IEEE.

Jingqi Feng, Yukai Huang, Rui Zhang, Sicheng Liang,  
Ming Yan, and Jie Wu. 2025. Windserv: Efficient  
phase-disaggregated llm serving with stream-based  
dynamic scheduling. In *Proceedings of the 52nd  
Annual International Symposium on Computer Archi-  
tecture*, pages 1283–1295.

Juncheng Gu, Youngmoon Lee, Yiwen Zhang, Mosharaf  
Chowdhury, and Kang G. Shin. 2017. *Efficient mem-  
ory disaggregation with infiniswap*. In *14th USENIX  
Symposium on Networked Systems Design and Imple-  
mentation (NSDI 17)*, pages 649–667, Boston, MA.  
USENIX Association.

Guseul Heo, Sangyeop Lee, Jaehong Cho, Hyunmin  
Choi, Sanghyeon Lee, Hyungkyu Ham, Gwang-  
sun Kim, Divya Mahajan, and Jongse Park. 2024.  
Neupims: Npu-pim heterogeneous acceleration for  
batched llm inferencing. In *Proceedings of the 29th  
ACM International Conference on Architectural Sup-  
port for Programming Languages and Operating Sys-  
tems, Volume 3*, pages 722–737.

Ke Hong, Lufang Chen, Zhong Wang, Xiuhong Li, Qi-  
uli Mao, Jianping Ma, Chao Xiong, Guanyu Wu,  
Buhe Han, Guohao Dai, and 1 others. 2025. semi-  
pd: Towards efficient llm serving via phase-wise dis-  
aggregated computation and unified storage. *arXiv  
preprint arXiv:2504.19867*.

Ke Hong, Guohao Dai, Jiaming Xu, Qiuli Mao, Xi-  
uhong Li, Jun Liu, Kangdi Chen, Yuhang Dong, and  
Yu Wang. 2024a. Flashdecoding++: Faster large lan-  
guage model inference with asynchronization, flat  
gemma optimization, and heuristics. *Proceedings of  
Machine Learning and Systems*, 6:148–161.

Ke Hong and 1 others. 2024b. Flashdecoding++: Faster  
large language model inference with asynchronization,  
flat gemm optimization, and heuristics. *Pro-  
ceedings of Machine Learning and Systems (MLSys)*,  
6:148–161.

Cunchen Hu, Heyang Huang, Junhao Hu, Jiang Xu,  
Xusheng Chen, Tao Xie, Chenxi Wang, Sa Wang,  
Yungang Bao, Ninghui Sun, and 1 others. 2024a.  
Memserv: Context caching for disaggregated llm  
serving with elastic memory pool. *arXiv preprint  
arXiv:2406.17565*.

Junda Hu, Yinmin Zhong, Shengyu Liu, Jianbo Chen,  
Xuanzhe Liu, and Hao Zhang. 2024b. TetriInfer:  
Disaggregated locality-aware inference for LLMs.  
In *Proceedings of the 29th ACM International Con-  
ference on Architectural Support for Programming  
Languages and Operating Systems (ASPLOS)*.

Jiayu Huang, Hailong Yang, Anthony Cheung, Shan  
Lu, Michael Leach, Aastha Mehta, Ayan Sinha, Prab-  
hakar Misra, Karthikeyan Shanmugam, and Pradeep

662	Dubey. 2022. <a href="#">Pond: CXL-based memory pooling systems for cloud platforms</a> . In <i>Proceedings of the ACM Symposium on Cloud Computing (SoCC '22)</i> .	Foteini Strati, Sara Mcallister, Amar Phanishayee, Jakob Tarnawski, and Ana Klimovic. 2024. <a href="#">Déjàvu: Kvcache streaming for fast, fault-tolerant generative llm serving</a> . <i>Preprint</i> , arXiv:2403.01876.	718
663			719
664			720
665	Siddharth Jha, Coleman Hooper, Xiaoxuan Liu, Sehoon Kim, and Kurt Keutzer. 2024. Learned best-effort llm serving. <i>arXiv preprint arXiv:2401.07886</i> .	Chao Wang, Pengfei Zuo, Zhangyu Chen, Yunkai Liang, Zhou Yu, and Ming-Chang Yang. 2025a. Prefill-decode aggregation or disaggregation? unifying both for goodput-optimized llm serving. <i>arXiv preprint arXiv:2508.01989</i> .	722
666			723
667			724
668	Youhe Jiang, Ran Yan, and Binhang Yuan. 2025. <a href="#">Hexgen-2: Disaggregated generative inference of llms in heterogeneous environment</a> . <i>Preprint</i> , arXiv:2502.07903.		725
669			726
670			727
671		Yuxin Wang, Yuhan Chen, Zeyu Li, Xueze Kang, Yuchu Fang, Yeju Zhou, Yang Zheng, Zhenheng Tang, Xin He, Rui Guo, and 1 others. 2025b. <a href="#">Burstgpt: A real-world workload dataset to optimize llm serving systems</a> . In <i>Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2</i> , pages 5831–5841.	728
672	Yibo Jin, Tao Wang, Huimin Lin, Mingyang Song, Peiyang Li, Yipeng Ma, Yicheng Shan, Zhengfan Yuan, Cailong Li, Yajing Sun, and 1 others. 2024. <a href="#">P/d-serve: Serving disaggregated large language model at scale</a> . <i>arXiv preprint arXiv:2408.08147</i> .		729
673			730
674			731
675			732
676			733
677	Aditya K Kamath, Ramya Prabhu, Jayashree Mohan, Simon Peter, Ramachandran Ramjee, and Ashish Panwar. 2025. <a href="#">Pod-attention: Unlocking full prefill-decode overlap for faster llm inference</a> . In <i>Proceedings of the 30th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2</i> , pages 897–912.	Juntao Zhao, Borui Wan, Yanghua Peng, Haibin Lin, and Chuan Wu. 2024. <a href="#">Llm-pq: Serving llm on heterogeneous clusters with phase-aware partition and adaptive quantization</a> . <i>arXiv preprint arXiv:2403.01136</i> .	734
678			735
679			736
680			737
681		Yinmin Zhong, Shengyu Liu, Jianbo Chen, Junda Hu, Yibo Zhu, Xuanzhe Liu, Xin Jin, and Hao Zhang. 2024. <a href="#">DistServe: Disaggregating pre-fill and decoding for goodput-optimized large language model serving</a> . In <i>18th USENIX Symposium on Operating Systems Design and Implementation (OSDI 24)</i> .	738
682			739
683			740
684	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. <a href="#">Efficient memory management for large language model serving with pagedattention</a> . In <i>Proceedings of the 29th Symposium on Operating Systems Principles (SOSP '23)</i> .		741
685			742
686			743
687			744
688			745
689			746
690			747
691	Kay Ousterhout, Aurojit Panda, Joshua Rosen, Ion Morgenthaler, Simon andMatin, Scott Shenker, and Sylvia Ratnasamy. 2013. <a href="#">Sparrow: Distributed, low latency scheduling</a> . In <i>Proceedings of the 24th ACM Symposium on Operating Systems Principles (SOSP '13)</i> , pages 69–84.	Zixuan Zhou, Xuefei Ning, Ke Hong, Tianyu Fu, Jiaming Xu, Shiyao Li, Yuming Lou, Luning Wang, Zhihang Yuan, Xiuhong Li, and 1 others. 2024. <a href="#">A survey on efficient inference for large language models</a> . <i>arXiv preprint arXiv:2404.14294</i> .	748
692			749
693			750
694			751
695			752
696		Jidong Zhuang and 1 others. 2024. <a href="#">Vserve: A context-aware serving system for large language models</a> . In <i>Proceedings of the 2024 USENIX Annual Technical Conference (ATC 24)</i> .	751
697	Pratyush Patel, Esha Choukse, Chaojie Zhang, Aashaka Shah, Íñigo Goiri, Saeed Maleki, and Ricardo Bianchini. 2024. <a href="#">Splitwise: Efficient generative llm inference using phase splitting</a> . In <i>2024 ACM/IEEE 51st Annual International Symposium on Computer Architecture (ISCA)</i> , pages 118–132. IEEE.		752
698			
699			
700			
701			
702			
703	Yuxin Qin, Weqning Li, , and 1 others. 2024. <a href="#">Mooncake: A kvcache-centric disaggregated architecture for llm serving</a> . <i>arXiv preprint arXiv:2407.00000</i> .		
704			
705			
706	Malte Schwarzkopf, Andy Konwinski, Michael Abd-El-Malek, and John Wilkes. 2013. <a href="#">Omega: Flexible, scalable schedulers for large compute clusters</a> . In <i>Proceedings of the 8th ACM European Conference on Computer Systems (EuroSys '13)</i> .		
707			
708			
709			
710			
711	Noam Shazeer. 2019. <a href="#">Fast transformer decoding: One write-head is all you need</a> . <i>arXiv preprint arXiv:1911.02150</i> .		
712			
713			
714	Xiaoxiang Shi, Colin Cai, Junjia Du, Zhanda Zhu, and Zhihao Jia. 2025. <a href="#">Nexus: Taming throughput-latency tradeoff in llm serving via efficient gpu sharing</a> . <i>arXiv e-prints</i> , pages arXiv–2507.		
715			
716			
717			

## A Additional Evidence for Prefill Queue Divergence

As shown in Figure 7, under memory-aware scheduling, *prefill* requests accumulate catastrophically on “Phantom” nodes—Node0’s queue reaches 40+ requests (peaking at 46), while other prefill workers maintain queues of only  $\sim 19$ . In contrast, compute-aware scheduling distributes prefill load evenly across all prefill workers ( $\sim 16$  requests each). This represents a  $2.4\times$  load skew, with Node0’s queue growing  $4.6\times$  from baseline while queues on other prefill workers merely double.

Figure 7(b) quantifies this prefill-phase disparity: during the burst plateau (300–420s), Node0 accumulates an average queue length of 36.2 requests under memory-aware scheduling, compared to 15.1 for other prefill workers and a uniform  $\sim 15.8$  under compute-aware scheduling.

## B Parameter Calibration and Configuration

To ensure reproducibility, this appendix details how Phasor calibrates the time-bound cost model parameters and configures control thresholds in our evaluation. Unless otherwise specified, all values refer to our V100 (FP16) serving stack.

**Notation.** For a prefill worker  $w \in \mathcal{P}$ , the time-bound virtual load aggregates per-request costs in its queue  $\mathcal{Q}_w$ :  $\mathcal{U}_{\mathcal{T}}(w) = \sum_{r \in \mathcal{Q}_w} (\alpha l_{in}(r)^2 + \beta P_m l_{in}(r))$ . We interpret  $\mathcal{U}_{\mathcal{T}}$  as a *FLOP-equivalent* proxy, and normalize it by the compute budget  $\mathcal{C}_{\mathcal{T}} = \mathbb{P}_{peak} T_{SLO}$  so that  $\mathcal{U}_{\mathcal{T}}/\mathcal{C}_{\mathcal{T}}$  approximates the fraction of the latency budget consumed by the queued work.

### B.1 Peak Throughput Measurement ( $\mathbb{P}_{peak}$ )

Theoretical peak numbers can deviate from the effective throughput of an inference stack. We therefore use  $\mathbb{P}_{peak}$  as an *effective peak throughput* measured on the target GPU under the same precision and kernel configuration as serving. Concretely, we run a GEMM-heavy microbenchmark representative of the prefill kernels and take the maximum sustained FP16 throughput as  $\mathbb{P}_{peak}$ . In our V100 setup, we observe  $\mathbb{P}_{peak} \approx 121$  TFLOP (FP16, tensor-core enabled). We report the measured effective peak rather than vendor-only specifications to keep  $\mathcal{U}_{\mathcal{T}}/\mathbb{P}_{peak}$  consistent with observed latency.

### B.2 Coefficient Profiling ( $\alpha, \beta$ )

The coefficients map request lengths to FLOP-equivalent costs under our serving stack.

**Reference model and normalization.** We use Qwen2.5-7B-Instruct (FP16) as the reference model  $M_{ref}$  and normalize the model-scale factor to it.

**Measurement.** We measure prefill latency  $T_{prefill}(l_{in})$  for synthetic requests with  $l_{in} \in [128, 8192]$  under a fixed batching policy to remove scheduler-induced variance.

**FLOP-equivalent conversion and fitting.** While  $\alpha$  and  $\beta$  could be derived from analytical FLOPs, such derivations are highly sensitive to implementation details (e.g., fused kernels, attention variants, and batching) and may deviate from the latency behavior observed in an end-to-end serving stack. Since scheduling decisions ultimately aim to control user-visible latency (e.g., TTFT) and prevent congestion under bursts, we calibrate the time-bound model using the most directly relevant observable: prefill latency. We therefore convert measured prefill latency into a FLOP-equivalent target:  $U(l_{in}) \triangleq T_{prefill}(l_{in}) \cdot \mathbb{P}_{peak}$ . We then fit  $U(l_{in}) \approx \alpha l_{in}^2 + \beta P_m(M_{ref}) l_{in}$  with least-squares regression ( $R^2 > 0.99$ ) to obtain stack-specific  $(\alpha, \beta)$ .

The fitted values used in our experiments are listed in Table 4.

### B.3 Cross-Model Generalization of the Time-Bound Cost Model

To improve cross-model generalization without per-model profiling of the linear component, we introduce an explicit model-scale factor  $P_m$  extracted from static metadata.

**Definition.** For dense Transformer LLMs, we define

$$P_m(M) \triangleq \frac{N_{\text{param}}(M)}{N_{\text{param}}(M_{ref})}, \quad (3)$$

where  $N_{\text{param}}(M)$  denotes the number of trainable parameters of model  $M$  obtained from model metadata. Thus, for our two evaluated models,  $P_m$  differs only by the parameter-scale ratio (e.g., 7B vs. 8B).

**Scope and boundary.** The above generalization targets *model changes under a fixed serving stack*. If the execution stack changes materially (e.g., precision, attention kernel, compiler/runtime),  $\mathbb{P}_{peak}$  and/or  $(\alpha, \beta)$  should be re-measured following §B.1–§B.2.

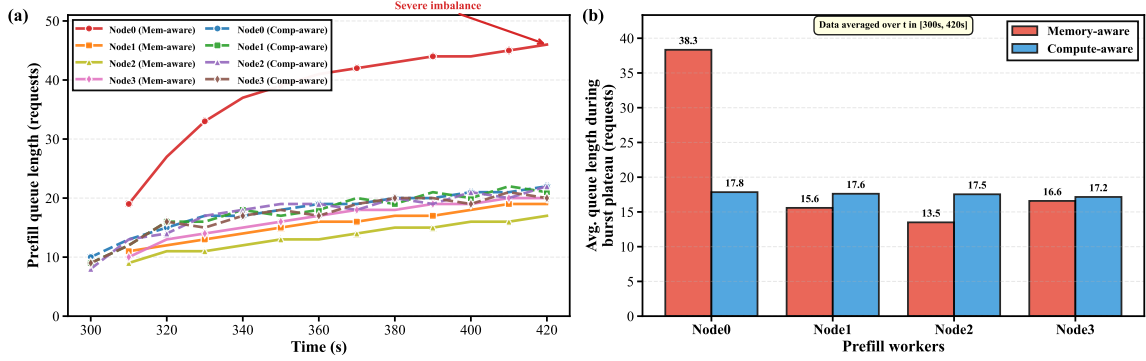


Figure 7: **Queue Divergence under Metric Misalignment.** Memory-aware routing piles requests on compute-saturated prefill workers, creating stragglers and amplifying head-of-line blocking; compute-aware routing maintains near-uniform queues.

#### B.4 Handling Negative Freeness

It is important to note that  $\Phi(w)$  can become negative if the estimated virtual load exceeds the physical budget (i.e.,  $\mathcal{U}_{virt} > \mathcal{C}_{phys}$ ). Physically, a negative  $\Phi(w) < 0$  signifies that the worker is currently overloaded relative to the strict SLO target.

**Graceful Degradation:** In such saturation regimes, the argmax operator in our Micro-scale Routing (Algorithm 1) remains valid and robust. By selecting the worker with the maximum algebraic value (i.e., the least negative  $\Phi$ ), Phasor naturally dispatches requests to the worker with the smallest overload margin.

**Interpretation:** This ensures that under extreme bursts where all workers violate SLOs, the system does not reject requests or deadlock; instead, it optimally distributes the excess load to minimize the aggregate magnitude of SLO violations, achieving graceful degradation.

#### B.5 Control Threshold Derivation

$$(\theta_{frag}, \theta_{flow})$$

Thresholds are chosen via sensitivity analysis to balance benefit and actuation overhead.

**De-fragmentation trigger threshold ( $\theta_{frag}$ ).** We sweep  $\theta_{frag} \in [0.05, 0.30]$  on a decode-heavy workload from sharegpt and evaluate (i) OOM/allocation-failure incidence and (ii) background maintenance overhead (extra time/bandwidth consumed by ASYNCDEFRAG). We select  $\theta_{frag} = 0.168$  as a Pareto point that eliminates OOM events in our testbed while keeping defragmentation overhead below a small budget (under 2% in our measurements).

**Flow control ratio ( $\theta_{flow}$ ).**  $\theta_{flow}$  provides hysteresis for macro-scale pool reconfiguration. We sweep

$\theta_{flow}$  and measure (i) reconfiguration frequency (to detect flapping) and (ii) end-to-end capacity/latency improvement under drifting workloads. We choose  $\theta_{flow} = 0.62$  as the smallest ratio that suppresses oscillatory reconfigurations while preserving most of the elasticity gains in our evaluation. In addition to this spatial hysteresis, we enforce a temporal Cool-down Period ( $T_{cool}$ ) to strictly prevent rapid role flipping. A worker promoted from the Mixed Pool is locked in its new role for a minimum duration (set to  $20 \times \Delta t$ , i.e., 1 second) before it can be recalled. This dual-mechanism ensures that reconfiguration only responds to sustained load shifts rather than transient noise.

#### B.6 Control Epoch ( $\Delta t$ )

$\Delta t$  sets the cadence of meso-/macro-scale control (e.g., de-fragmentation checks and pool rebalancing). It trades off responsiveness and control overhead: smaller  $\Delta t$  reacts faster but increases telemetry and background actuation frequency, while larger  $\Delta t$  reduces overhead but delays adaptation under bursts and workload drifts. We set  $\Delta t = 50$  ms in our testbed, matching the natural meso-scale timescale of decode-side memory dynamics while remaining orders of magnitude larger than micro-scale request routing. We report a sensitivity sweep of  $\Delta t$  in Appendix F.

#### B.7 Latency Budget ( $T_{SLO}$ )

We adopt the standard two-metric SLO formulation in LLM serving, where a request is considered *SLO-satisfied* only if it meets *both* a TTFT budget and a TPOT budget. In contrast, our time-bound freeness for *prefill* normalizes only by  $T_{SLO}^{TTFT}$ , since TTFT is determined by the prefill path, while TPOT is governed by the decode path.

Table 3: **Per-workload SLO budgets used in evaluation.** A request is counted as SLO-satisfied only if it meets *both* TTFT and TPOT budgets.

Workload ID	$T_{SLO}^{TTFT}$ (s)	$T_{SLO}^{TPOT}$ (s)
ShareGPT-S-S (Short QA)	0.30	0.15
ShareGPT-M-M (Dialog)	0.40	0.20
BurstGPT	0.40	0.20
ShareGPT-S-L (Code Gen.)	0.125	0.20
ShareGPT-L-S (Summarization)	2.25	0.13
ShareGPT-L-L (Translation)	3.00	0.18

SLO budgets are operator-configurable and workload-/tier-dependent in practice. Rather than claiming universal budgets, we align our per-scenario settings with widely-used tiers reported by prior serving systems (Zhong et al., 2024; Hong et al., 2025) and evaluations. Specifically, we map our synthetic archetypes to the closest prior workload categories (interactive chat, long-context reading, code completion) and use their published TTFT/TPOT targets as defaults. Table 3 summarizes the budgets used in our evaluation.

## C Visualization of Control Mechanisms

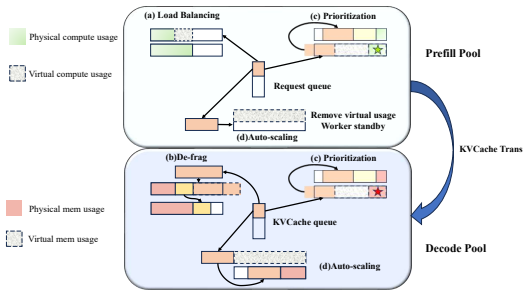


Figure 8: **Phasor Control Mechanisms.** Driven by Virtual Freeness, the system executes (a) Gradient-based Routing to balance compute, (b) Proactive De-fragmentation to reclaim memory, and (c/d) Prioritization and Auto-scaling to manage capacity.

## D Zero-Downtime KV Migration

Phasor introduces a multi-stage pipelining mechanism for KV migration. Instead of pausing execution, the system asynchronously replicates the generated KV Cache from the source to the target worker while the source continues inference. This approach effectively masks the data transfer latency, requiring only a negligible interruption for the final consistency check.

Specifically, to mitigate rescheduling overhead, Phasor employs an asynchronous, iterative migration strategy that pipelines state transfer with active decoding. As illustrated in Figure 9, the process exploits the append-only nature of the KV cache to minimize the critical section (downtime).

**Initial Copy Phase (Stage 0):** Upon triggering migration, the Source Worker continues decoding new tokens while asynchronously transferring the existing snapshot of KV blocks to the Destination Worker. This ensures the request makes progress during the bulk transfer.

**Page Iteration (Stage 1...N):** Since decoding continues during the initial phase, new "dirty" KV blocks are generated. Phasor enters an iterative convergence loop where it transfers only the incremental KV blocks generated in the previous stage. Because transfer speed typically exceeds generation speed, the volume of dirty data shrinks with each iteration.

**Stop-and-Copy Phase:** Once the remaining delta is small enough (Stage N), Phasor suspends execution for a final, atomic synchronization. This transfers the last few blocks and hands over control to the destination worker.

This design ensures that the Downtime is constant and independent of the total sequence length, limited only to the transfer time of the final iteration's output.

## E Comprehensive Performance Panorama

While the main text focuses on extreme bottleneck regimes, Figure 10 illustrates our holistic evaluation across the full spectrum of workload archetypes, including Quick Q&A (Short-Short), Translation (Long-Long), and Dialog (Medium-Medium).

### Analysis of Auxiliary Workloads.

1. Quick Q&A (Sharegpt-S-S): Low Overhead Validation.

**Data Overview:** This lightweight scenario ( $L_{in}, L_{out} < 512$ ) serves as a litmus test for system overhead. Across all load levels ( $\lambda = 5$  to 40 req/s), Phasor reduces the average Total Request P99 by 11.5% compared to Splitwise and 8.9% compared to DistServe.

**Highlight:** At peak load ( $\lambda = 40$ ), Phasor achieves a Total Request P99 of 28.9 s, outperforming Splitwise (33.5 s) by 13.8% and DistServe (32.4 s) by 10.9%. The TTFT P99 is reduced

Table 4: Phasor Parameters in Our V100 (FP16) Testbed

Symbol	Definition	Value
<i>Physical / Stack Constants</i>		
$\mathbb{P}_{peak}$	Measured effective peak throughput (FP16, serving stack)	$\approx 121$ TFLOPS
$\mathbb{M}_{total}$	Physical memory limit per decode worker	32 GB
$\alpha$	Quadratic cost coefficient (attention proxy, FLOP-equivalent)	$4.25 \times 10^{-5}$
$\beta$	Linear cost coefficient (FFN proxy, FLOP-equivalent; scaled by $P_m$ )	$6.80 \times 10^{-3}$
$P_m$	Model-scale factor	model-dependent
<i>Control Parameters</i>		
$\Delta t$	Control epoch / sampling interval	50 ms
$\theta_{frag}$	De-fragmentation trigger threshold	0.168
$\theta_{flow}$	Flow control ratio (hysteresis for pool reconfiguration)	0.62

990 to 202 ms vs Splitwise’s 215 ms and DistServe’s  
 991 208 ms, confirming minimal scheduling overhead  
 992 even for short tasks.

993 Attribution: These results contradict the intuition  
 994 that sophisticated schedulers penalize lightweight  
 995 requests. Phasor’s asynchronous control plane in-  
 996 curs negligible overhead ( $<3\%$  TTFT increase vs  
 997 optimal) while effectively preventing micro-burst  
 998 queuing through proactive load redistribution.

## 999 2. Translation (Sharegpt-L-L): Dual-Bottleneck 1000 Efficiency.

1001 Data Overview: In this heavy regime where both  
 1002 prefill and decode phases saturate ( $L_{in}, L_{out} \approx$   
 1003 2048), Phasor demonstrates an average Decode  
 1004 Time Per Token P99 reduction of 11.8% versus  
 1005 Splitwise and 4.1% versus DistServe across all ar-  
 1006 rival rates. For end-to-end performance, the aver-  
 1007 age Total Request P99 improves by 14.9% over  
 1008 Splitwise and 8.7% over DistServe.

1009 Highlight: At saturation load ( $\lambda = 40$ ), static  
 1010 partitioning collapses. DistServe degrades to a To-  
 1011 tal Request P99 of 175.9 s and Splitwise to 190.2 s,  
 1012 while Phasor suppresses this to 151.4 s—achieving  
 1013 a 20.4% reduction versus Splitwise and 14.0% ver-  
 1014 sus DistServe. The Decode Time Per Token P99 is  
 1015 maintained at 118.7 ms, significantly outperform-  
 1016 ing Splitwise (138.9 ms, +17.0% slower) and Dist-  
 1017 Serve (125.5 ms, +5.7% slower).

1018 Attribution: When both stages are bottlenecked,  
 1019 rigid resource boundaries become pathological.  
 1020 Phasor’s Mixed Pool enables elastic capacity real-  
 1021 location, dynamically borrowing idle decode GPUs  
 1022 for prefill bursts. This fluidity mitigates the “hard  
 1023 edges” of static splits, preventing cascading queue  
 1024 buildup when the entire cluster is under stress.

## 1025 3. Dialog (Sharegpt-M-M): Balanced Workload 1026 Robustness.

Data Overview: For standard balanced work- 1027  
 loads ( $L \approx 1024$ ), Phasor demonstrates the most 1028  
 pronounced superiority. Across all load levels, Pha- 1029  
 sor achieves an average Total Request P99 reduc- 1030  
 tion of 25.6% compared to Splitwise and 19.6% 1031  
 compared to DistServe. 1032

1033 Highlight: At peak load ( $\lambda = 40$ ), Phasor re- 1033  
 duces Total Request P99 by 30.7% compared to 1034  
 Splitwise (102.5 s vs 148.0 s) and 19.9% compared 1035  
 to DistServe (102.5 s vs 128.0 s). This advantage 1036  
 emerges without extreme workload skew—neither 1037  
 stage is exclusively bottlenecked, yet static sched- 1038  
 ulers still suffer from suboptimal partitioning under 1039  
 load variance. 1040

1041 Attribution: This confirms that Phasor acts as 1041  
 a robust generalist. Even when workloads lack 1042  
 extreme input/output length skew, production traf- 1043  
 fic exhibits inherent burstiness and heterogeneity. 1044  
 Phasor’s dimensionless interface unifies heteroge- 1045  
 neous resource demands (compute-bound prefill vs 1046  
 memory-bound decode) into a common optimiza- 1047  
 tion space, enabling gradient-based balancing that 1048  
 static heuristics cannot achieve. 1049

## 1050 F Sensitivity Analysis

1051 To evaluate the robustness of Phasor under varying 1051  
 hyper-parameter settings, we conduct a compre- 1052  
 hensive sensitivity analysis on the BurstGPT work- 1053  
 load. We examine four critical dimensions: the 1054  
 Cost Weight ( $\alpha$ ), the De-fragmentation Threshold 1055  
 ( $\theta_{frag}$ ), the Control Epoch ( $\Delta t$ ), and the Flow Con- 1056  
 trol Ratio ( $\theta_{flow}$ ). Table 5 summarizes the impact 1057  
 on P99 Latency, Capacity, and TTFT. 1058

1059 **Robustness Observations. Cost Weight ( $\alpha$ ):** 1059  
 The system exhibits asymmetric sensitivity. Under- 1060  
 estimating  $\alpha$  ( $0.5\times$ ) is detrimental, causing con- 1061

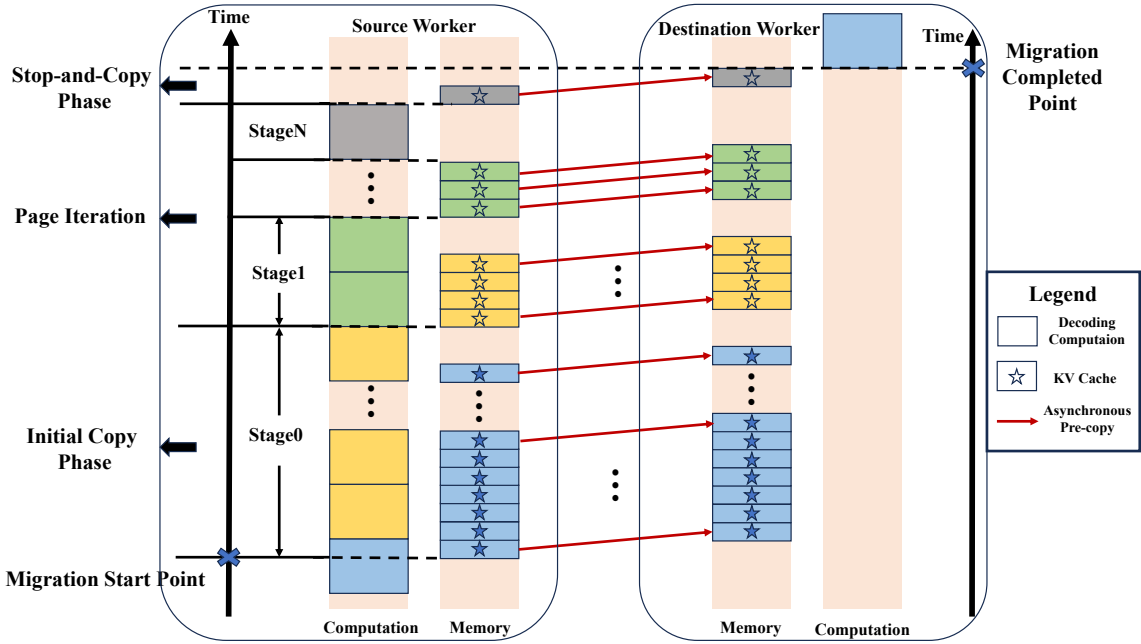


Figure 9: **Phasor’s Asynchronous Pipelined Migration Workflow.** The migration is divided into the *Initial Copy Phase* (bulk background transfer), *Page Iteration* (incremental state convergence), and a minimal *Stop-and-Copy Phase* (atomic handover). By pipelining KV transfer with decoding computation on the Source Worker, Phasor achieves near-zero downtime regardless of the context length.

1062 gestion and a 14% latency spike. Conversely, over-  
 1063 estimating  $\alpha$  ( $2.0\times$ ) makes the scheduler conserva-  
 1064 tive: it maintains strict latency compliance ( $P99 \approx$   
 1065  $0.99$ ) by actively rejecting risky batches. However,  
 1066 this safety comes at a cost, effectively "wasting"  
 1067 usable headroom and leading to a 6% drop in ser-  
 1068 viceable capacity.

1069 **De-fragmentation Threshold** ( $\theta_{frag}$ ): We ob-  
 1070 serve a saturation effect. While a low threshold  
 1071 ( $0.10$ ) allows fragmentation to degrade latency  
 1072 ( $+9\%$ ), raising the threshold beyond the default  
 1073 ( $0.35$ ) yields diminishing returns. It does not fur-  
 1074 ther improve latency (as blocking is already elim-  
 1075 inated) but incurs slightly higher migration over-  
 1076 head, reducing capacity by 2%.

1077 **Control Epoch** ( $\Delta t$ ): Phasor demonstrates remark-  
 1078 able robustness to the telemetry interval. Increas-  
 1079 ing  $\Delta t$  from 50ms to 200ms only degrades per-  
 1080 formance by 6–8%, designated as "Low" impact.  
 1081 This validates that the predictive nature of Virtual  
 1082 Freeness reduces the reliance on high-frequency  
 1083 hardware feedback.

1084 **Flow Control Ratio** ( $\theta_{flow}$ ): The system remains  
 1085 stable across a wide range of flow ratios. Because  
 1086  $\theta_{flow}$  works in tandem with the temporal Cool-  
 1087 down Period (Appendix B.4), the Mixed Pool re-  
 1088 sists oscillation even when the ratio is loose ( $0.50$ )  
 1089 or strict ( $0.80$ ), showing only marginal impact on

end-to-end metrics. 1090

## 1091 G Energy Efficiency Analysis 1092

1092 In the context of LLM efficiency, energy consump-  
 1093 tion is the physical denominator of utilization. We  
 1094 therefore test whether Phasor’s throughput gains  
 1095 translate into better *work per joule*, rather than  
 1096 merely higher power draw.

1097 **Methodology.** We monitor *aggregate GPU*  
 1098 *power* (sum over 8 GPUs) via `nvidia-smi/NVML`  
 1099 at 100 ms intervals and evaluate the *Production Mix*  
 1100 (*BurstGPT*) at the peak-load phase (40–60 req/s),  
 1101 excluding warm-up. We define **Energy Efficiency**  
 1102 as *Joules per generated output token* (J/token),  
 1103 computed as  $\overline{P_{GPU}} / \text{tok/s}$ , which is equivalent to  
 1104  $\int P(t) dt / \#\text{generated tokens}$  over the same win-  
 1105 dow. We use *output-token goodput* (success-  
 1106 fully generated output tokens/s) to reflect end-to-  
 1107 end generation progress, which is most sensitive  
 1108 to decode-side stalls and fragmentation-induced  
 1109 effective-capacity loss.

1110 **Sanity Check: Utilization vs. Power.** As shown  
 1111 in Table 6, Phasor draws higher average GPU  
 1112 power ( $\sim 1720$  W) than Splitwise ( $\sim 1315$  W), a  
 1113 31% increase. This is consistent with Phasor mit-  
 1114 igating pipeline starvation: by sustaining higher

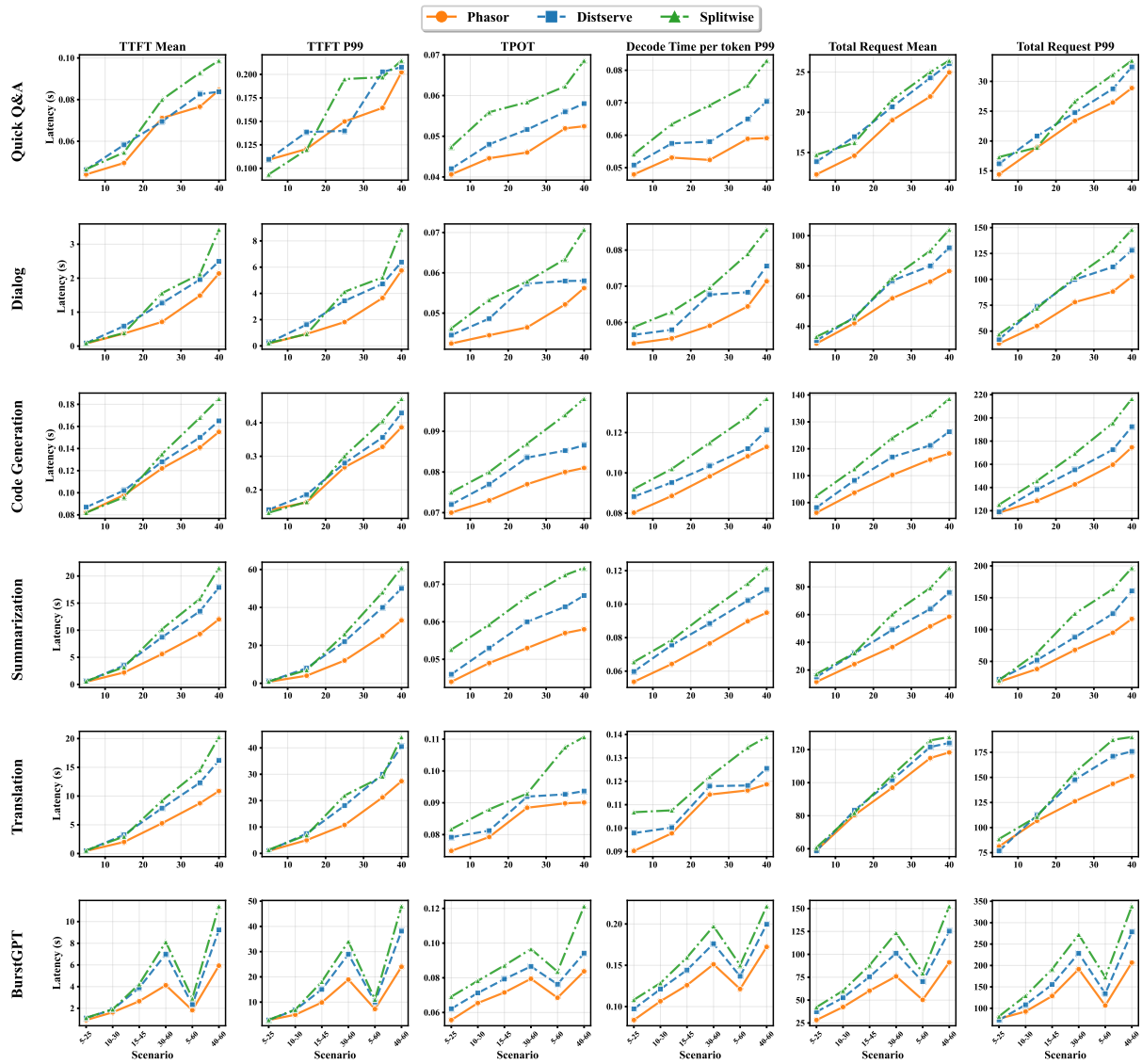


Figure 10: **Comprehensive Performance Panorama.** End-to-end latency metrics across six workload archetypes, covering Quick Q&A (S-S), Dialog (M-M), Code Generation (S-L), Summarization (L-S), Translation (L-L), and Production Mix (BurstGPT). Phasor demonstrates universal advantages in TTFT, TPOT, and total request latency across all load levels.

1115 effective decode concurrency, it converts idle GPU  
 1116 power into active generation work.

1117 **Conclusion.** Phasor improves energy proportion-  
 1118 ality by translating higher utilization into lower  
 1119 GPU energy per generated token, rather than  
 1120 merely increasing power draw.

Dimension	Parameter Variation	P99 Latency ( $\downarrow$ )	Capacity ( $\uparrow$ )	TTFT P99 ( $\downarrow$ )	Sensitivity
<b>Cost Weight</b> ( $\alpha$ )	$0.5 \times \alpha_{def}$	1.14	0.92	1.18	Moderate
	$2.0 \times \alpha_{def}$	<b>0.99</b>	0.94	1.00	Low
<b>De-frag Threshold</b> ( $\theta_{frag}$ )	Low (0.10)	1.09	0.96	1.12	Moderate
	High (0.35)	1.00	0.98	1.00	Low
<b>Control Epoch</b> ( $\Delta t$ )	Fast (20 ms)	0.99	1.01	0.99	Negligible
	Slow (200 ms)	1.06	0.97	1.08	Low
<b>Flow Ratio</b> ( $\theta_{flow}$ )	Loose (0.50)	1.02	0.99	1.01	Low
	Strict (0.80)	1.04	0.95	1.03	Low

Table 5: Sensitivity Analysis on BurstGPT. Metrics are normalized to the default configuration ( $\alpha_{def}, \theta_{frag} = 0.168$ ,  $\Delta t = 50\text{ms}$ ,  $\theta_{flow} = 0.62$ ). "Sensitivity" denotes the system's responsiveness to parameter drift.

System	Avg. Power (W)	Gen. Throughput (tok/s)	Efficiency (J/Token)
Splitwise	$\sim 1,315$	$\sim 5,422$	0.243
DistServe	$\sim 1,410$	$\sim 7,009$	0.201
<b>Phasor</b>	<b><math>\sim 1,720</math></b>	<b><math>\sim 9,820</math></b>	<b>0.175</b>

Table 6: Energy Efficiency on BurstGPT ( $8 \times V100$ ). Phasor converts "idle power" into "active work". Despite a higher raw power draw, the substantial gain in throughput results in a **28.0% reduction in energy cost per token**.