

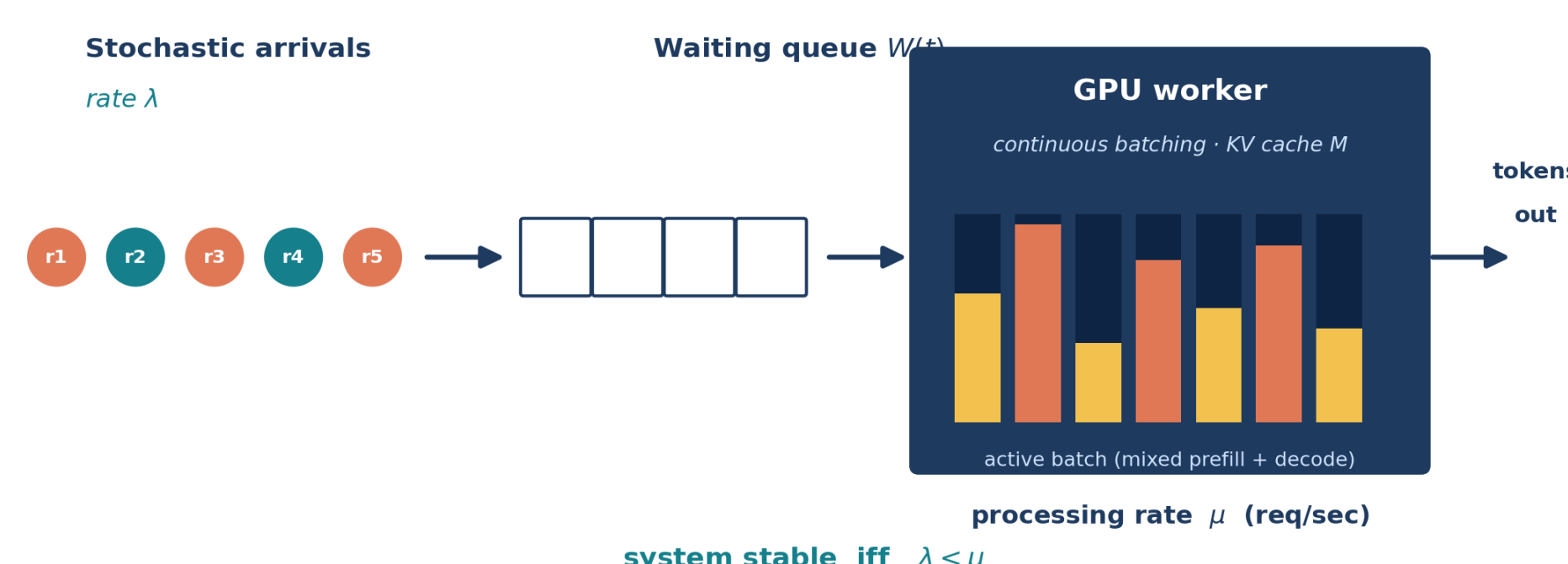
TL;DR

LLM inference at scale is bottlenecked by **BOTH compute AND KV-cache memory**. We build the **first queuing model that captures memory**, derive **closed-form stability conditions**, and validate on real A100 GPUs with GAP < 10%.

Central Question

Given a workload and GPU memory M , how many GPUs are needed to keep an LLM service stable?

LLM Inference as a Queue



Two phases per request: **prefill** (chunked, memory grows by \hat{s} each step) then **decode** (memory grows by 1 token per step). The GPU runs a mixed batch under a hard memory cap M .

Our Contributions

- Modeling**
First queuing model for LLM inference that explicitly captures both compute AND GPU memory (KV cache) constraints — generalizing prior compute-only formulations.
- Theory**
Rigorous stability & instability conditions via a Lyapunov / Foster argument adapted to the piecewise-linear KV growth.
- Validation**
Extensive real A100 experiments under single- and multi-GPU deployments, plus LongBench v2. Predictions agree to within 10%.

Why It Matters in Practice

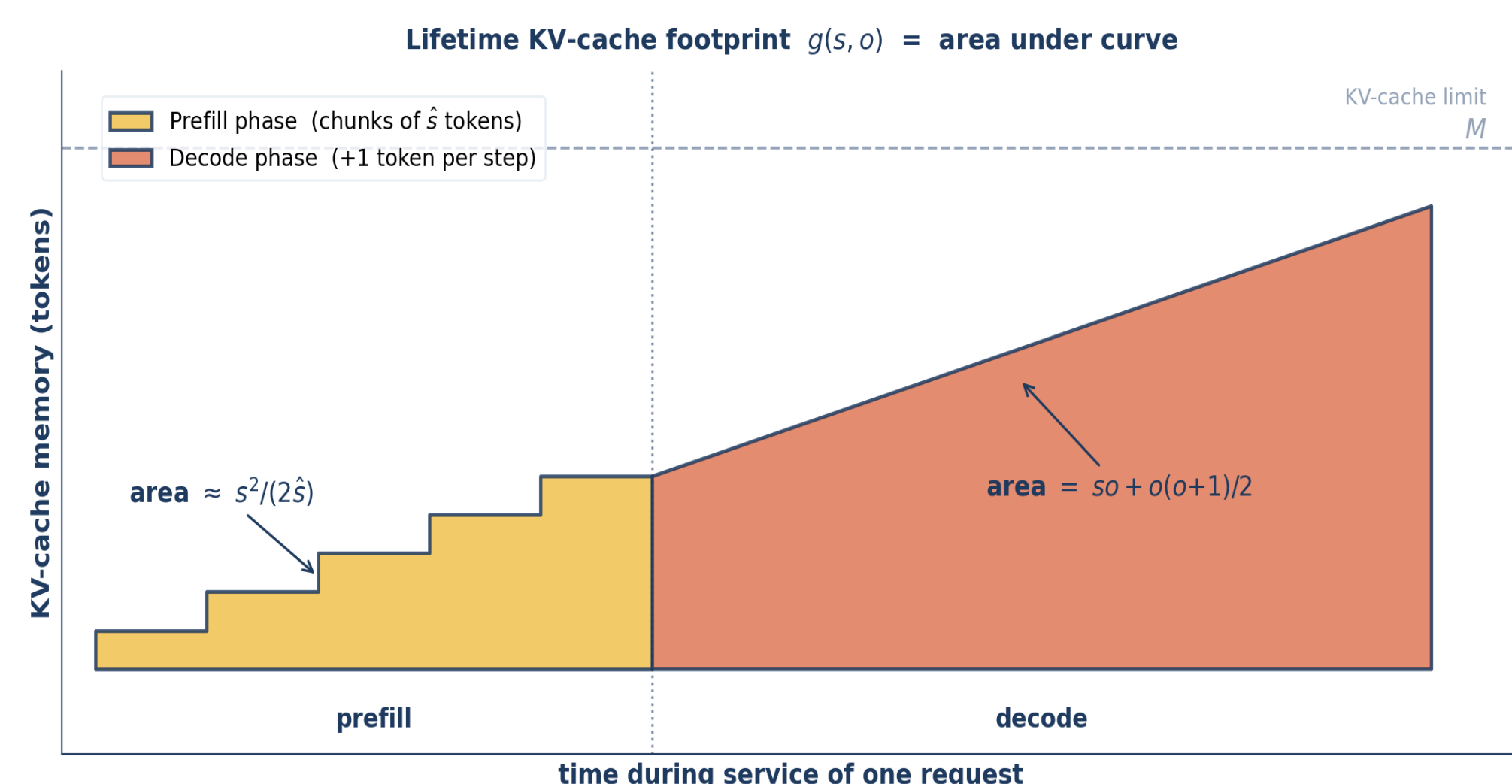
- Right-size clusters**
Pick $N = \lceil \lambda / \mu \rceil$ GPUs directly from the workload distribution — no costly trial-and-error rollouts.
- Avoid OOM at peak**
The memory-aware rate μ catches regimes that compute-only models silently misjudge as feasible.
- Adapt to drift**
Plug a fresh estimate of $p(s,o)$ into μ to autoscale as the workload mix shifts in production.

ICML 2026 · Full paper + proofs in supplementary

Memory Footprint of a Request

For input s , output o , and chunk size \hat{s} :

$$g(s, o) = \frac{(1 + s/\hat{s})s + 2so + o(o + 1)}{2}$$



Key insight: the lifetime area under the memory curve is what each request “costs” the GPU. Compressing the trajectory into the scalar $g(s, o)$ makes Lyapunov analysis tractable.

Why Constant Batch Time?

Near saturation, **attention compute dominates**; chunked prefill bounds per-batch work. Empirically **80% of batches share the same execution time** → we model \bar{b} as constant.

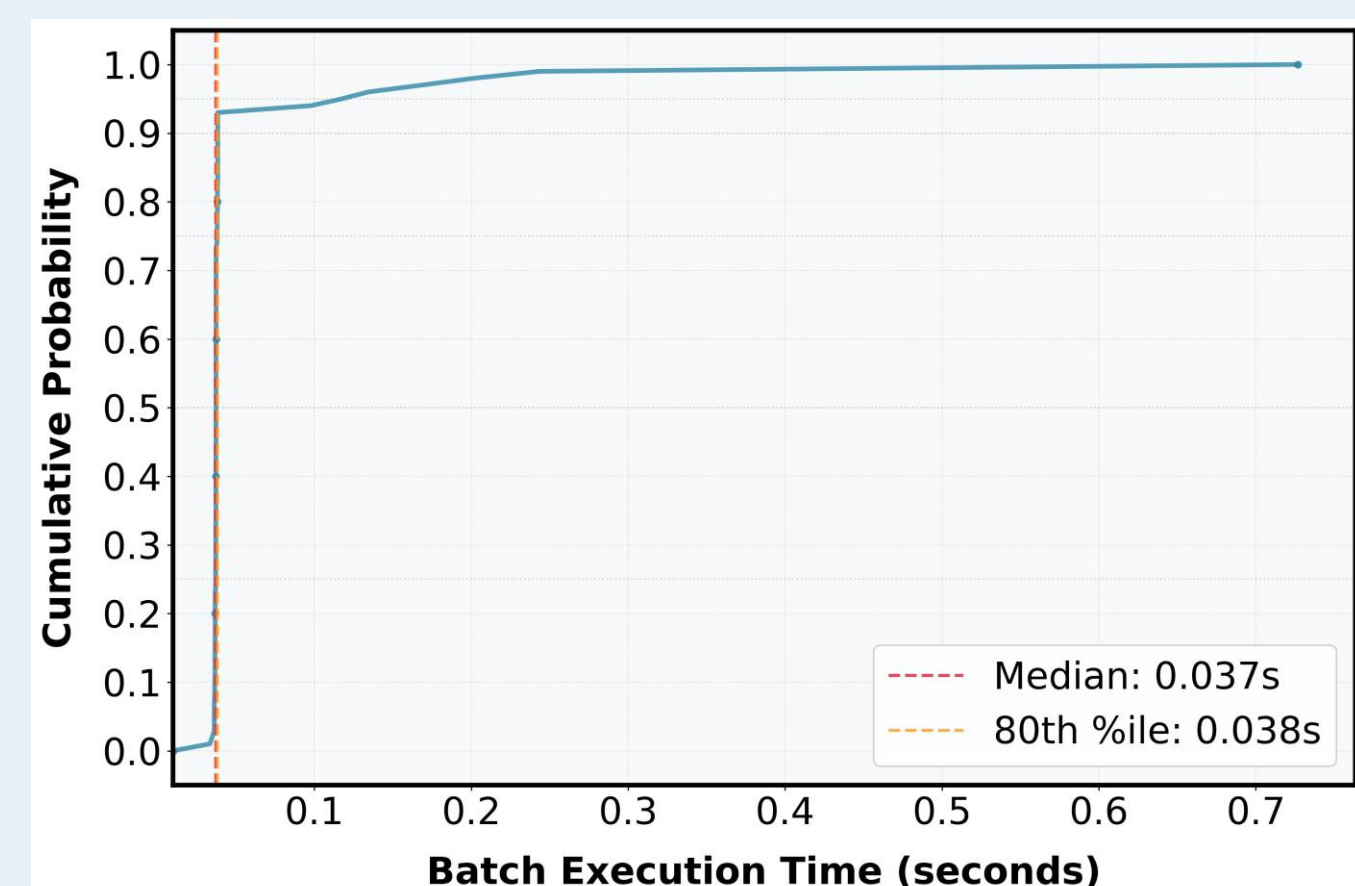
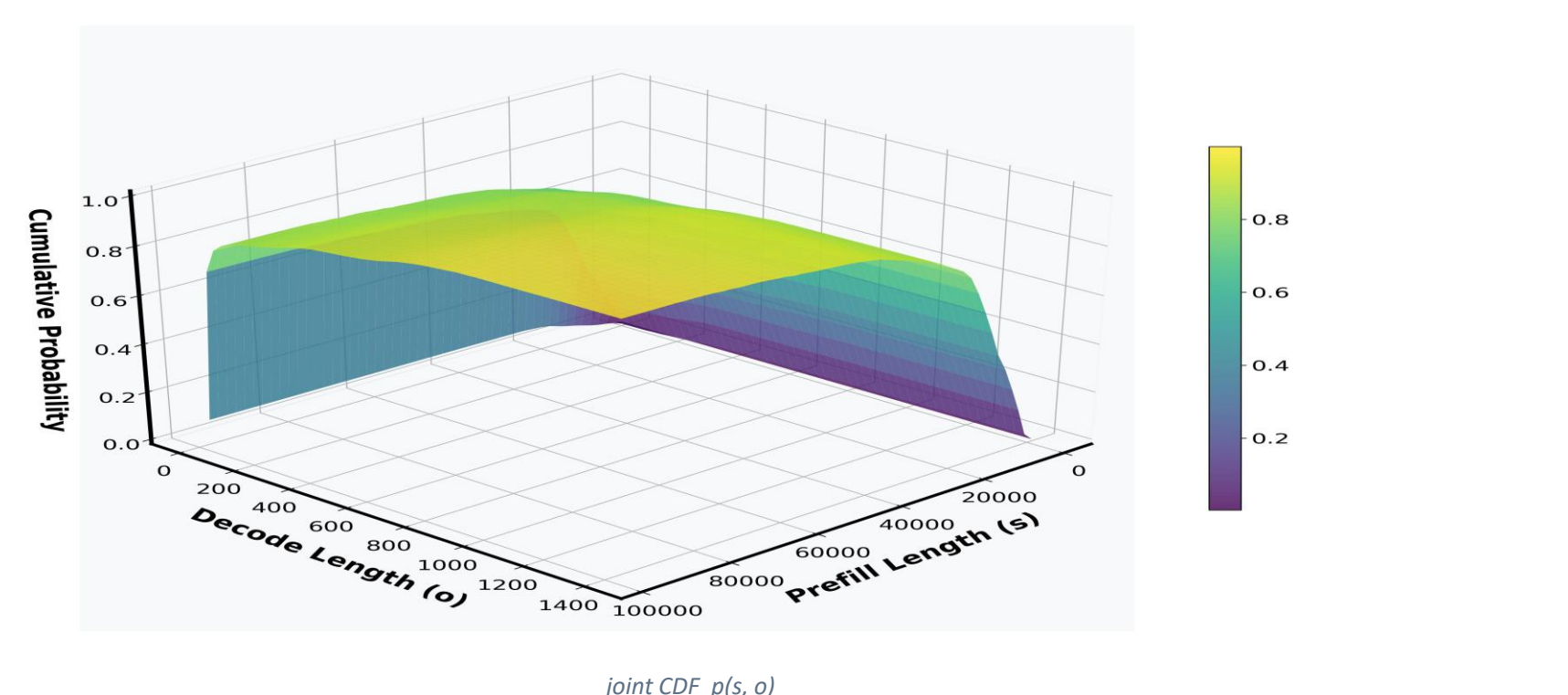
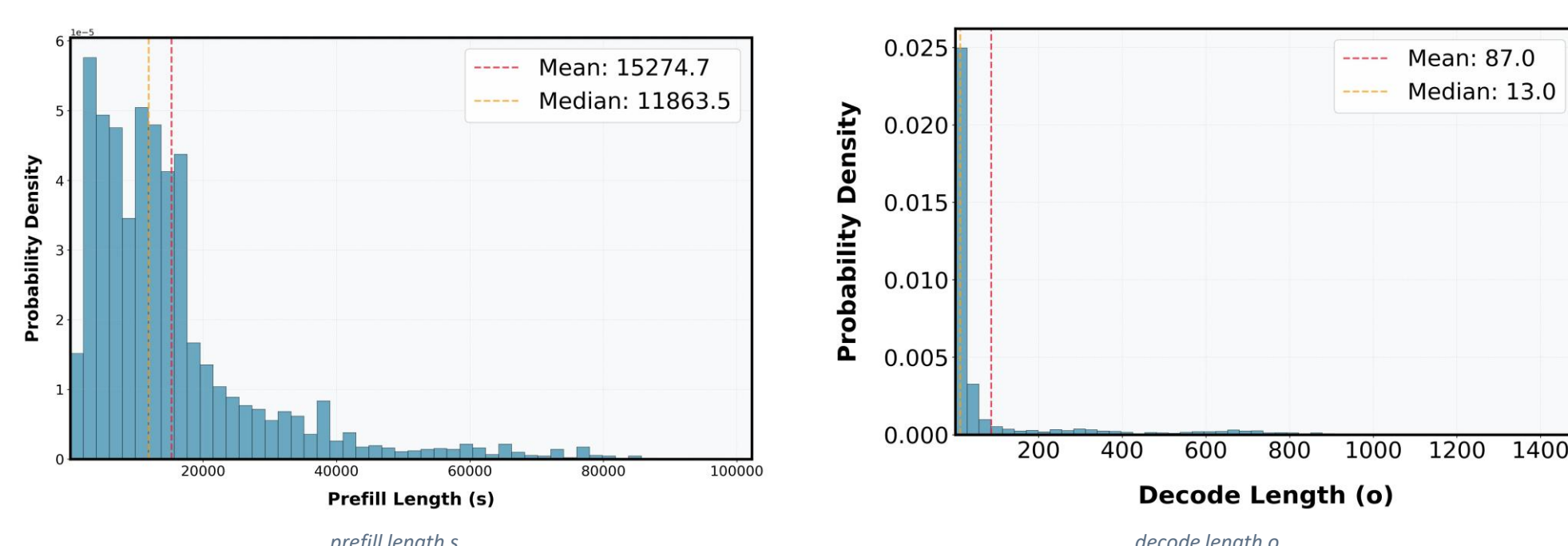


Fig. CDF of batch execution time (Llama-3-8B · A100 · PD 1:1). Median \approx 80th %ile

Heavy-Tailed Real Workloads

LongBench v2: 503 long-context queries — prefill is heavy-tailed, decode is heavier-tailed. Mean \gg median in both → modelling the full joint $p(s,o)$ is essential.



Main Result: Stability Condition

Per-GPU service rate

$$\mu = \frac{M}{\bar{b} \cdot \mathbb{E}[g(s, o) \sim p(s, o)]}$$

STABLE

$$\lambda < \mu(1 - \delta)$$

OVERLOAD

$$\lambda > \mu$$

THEOREM 1 · Overload

If $\lambda > \mu$, the queue grows to infinity under ANY scheduling policy.

THEOREM 2 · Stability

If $\lambda < \mu(1 - \delta)$, with $\delta = \text{esssup}(s+o)/M \ll 1$, any work-conserving + continuous-batching policy is positive recurrent.

Proof Sketch · Lyapunov Drift

Let $V(t)$ = total outstanding KV-cache demand (“area still to erase”). Per slot:

- Saturated GPU erases $\geq M(1 - \delta)$ per slot
- Arrivals add $\lambda \bar{b} \cdot \mathbb{E}[g(s, o)]$ in expectation

$$\mathbb{E}[V(t+1) - V(t) | \text{state}_t] \leq -M(1 - \delta) + \lambda \bar{b} \mathbb{E}[g(s, o)]$$

⇒ Negative drift outside a compact set when $\lambda < \mu(1 - \delta)$. Foster–Lyapunov ⇒ stability. □

GPU Provisioning Rule

$$N_{\text{GPU}} = \lceil \lambda / \mu \rceil$$

Once λ and μ are estimated, **right-size the cluster in one line**.

Worked Example · PD 1:1, Llama-3-8B on A100

GPU memory M 131 k tokens
Batch time \bar{b} 0.037 s

Theory μ_{theory} 3.263 req/s
Measured μ_{gpu} 3.387 req/s

Scope · Extensions · Future Work

- Tensor parallelism:** treat a TP group as one logical worker — formula for μ is unchanged, swap in TP-specific M and \bar{b} .
- Pipeline parallelism / PD disaggregation:** tandem queues with separate memory budgets — open theoretical direction.
- Speculative decoding & MoE:** draft acceptance rate and expert-routing skew enter $g(\cdot)$ multiplicatively — straightforward to plug in.
- Beyond stationarity:** use the per-slot drift bound as a building block for non-stationary $\lambda(t)$, online autoscaling, and SLA-aware admission control.

Experimental Validation

Testbed: 8 × NVIDIA A100, Meta-Llama-3-8B, vLLM v1, chunked prefill $\hat{s}=512$.
Memory: $M = 131$ k tokens. Batch time \bar{b} from median of 10k simulated batches.
Workloads: PD ratios 1:1, 2:1, 1:2, a piecewise-stationary mix, and LongBench v2.

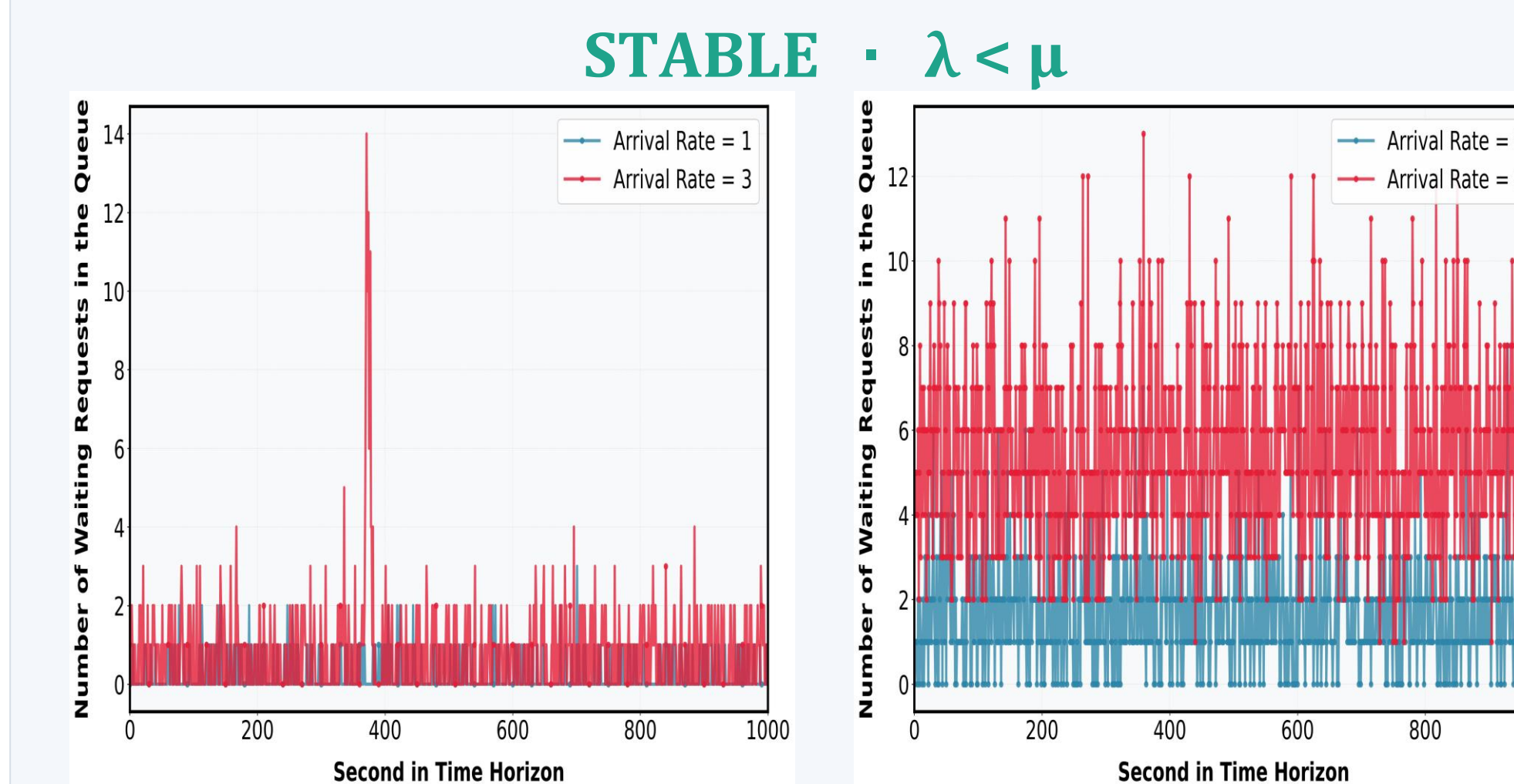
Theory Matches Reality (GAP < 10%)

Workload	μ_{gpu}	μ_{theory}	GAP
PD 1:1	3.387	3.263	3.66%
PD 2:1	3.650	3.956	8.38%
PD 1:2	2.969	2.902	2.25%
Mixed 2:1 → 1:2	3.137	3.385	7.90%
LongBench v2	0.610	0.561	8.03%
8 × GPU, PD 1:1	26.710	25.808	3.38%

Single- and multi-GPU predictions are within 10% of measurements across every regime — including a heavy-tailed real-world dataset.

Queue Dynamics (single & 8-GPU)

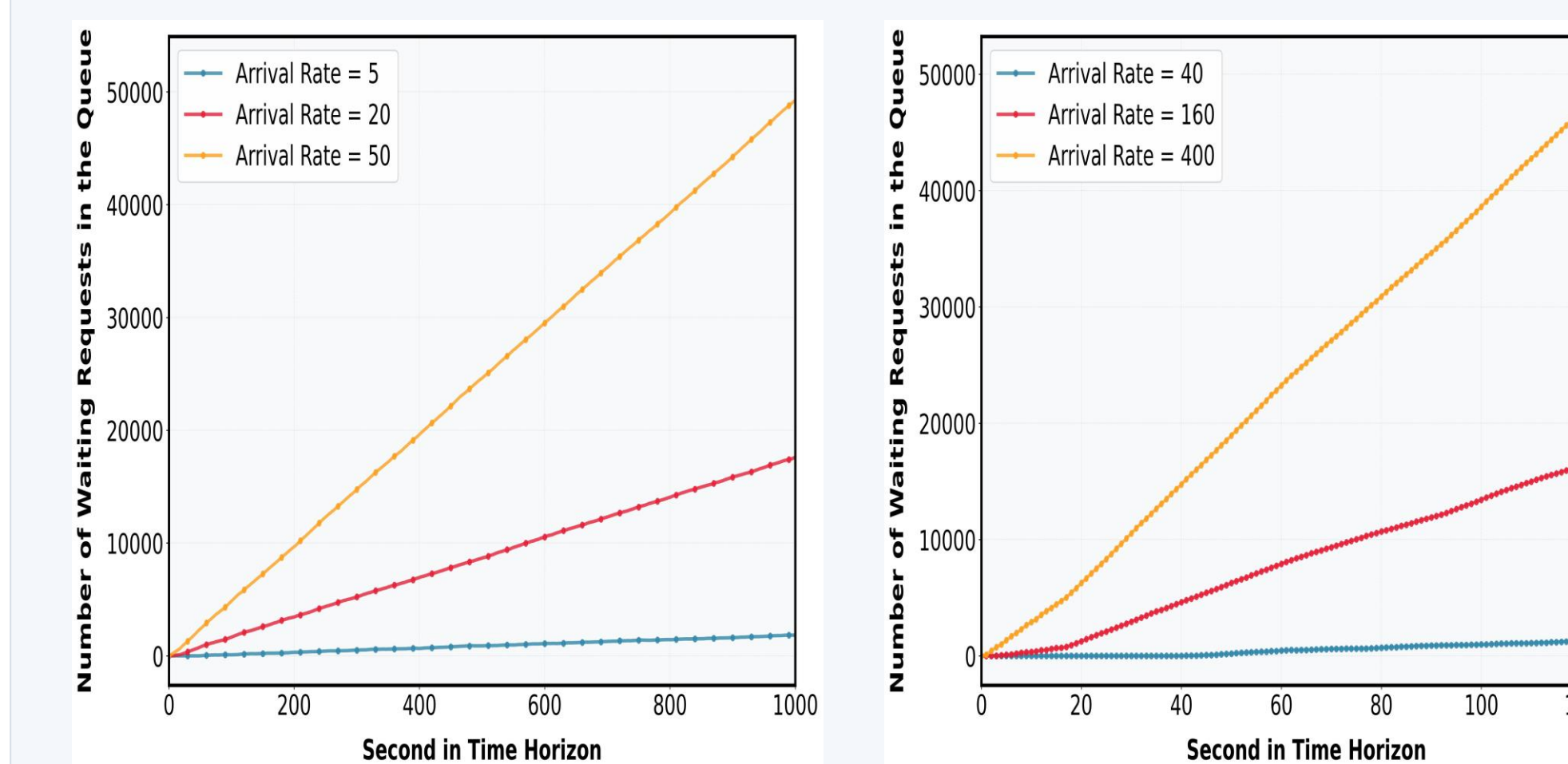
Same conclusion in both scales: when $\lambda < \mu$ the backlog stays bounded; when $\lambda > \mu$ it grows linearly in time.



1 GPU · queue stays bounded

8 GPUs · queue stays bounded

OVERLOAD · $\lambda > \mu$



1 GPU · linear blow-up

8 GPUs · linear blow-up