

# GPU Power Management Enables Rapid Deployment of Large Language Models

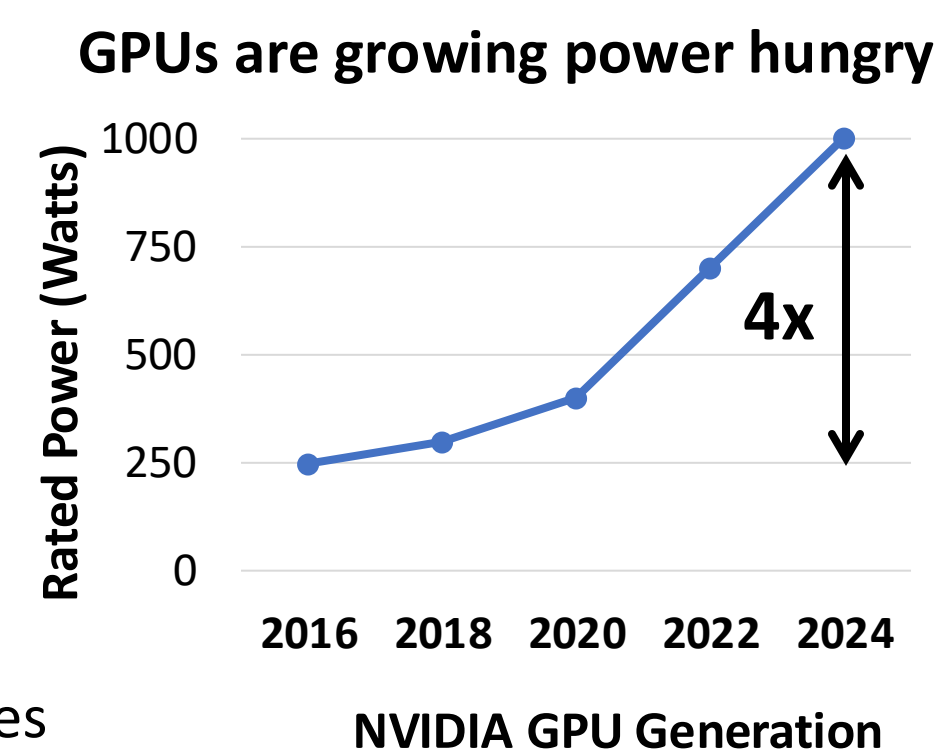
## Characterizing Power Management Opportunities for LLMs in the Cloud

Power is a key bottleneck for LLM deployments at scale

### Big Tech's Latest Obsession Is Finding Enough Energy

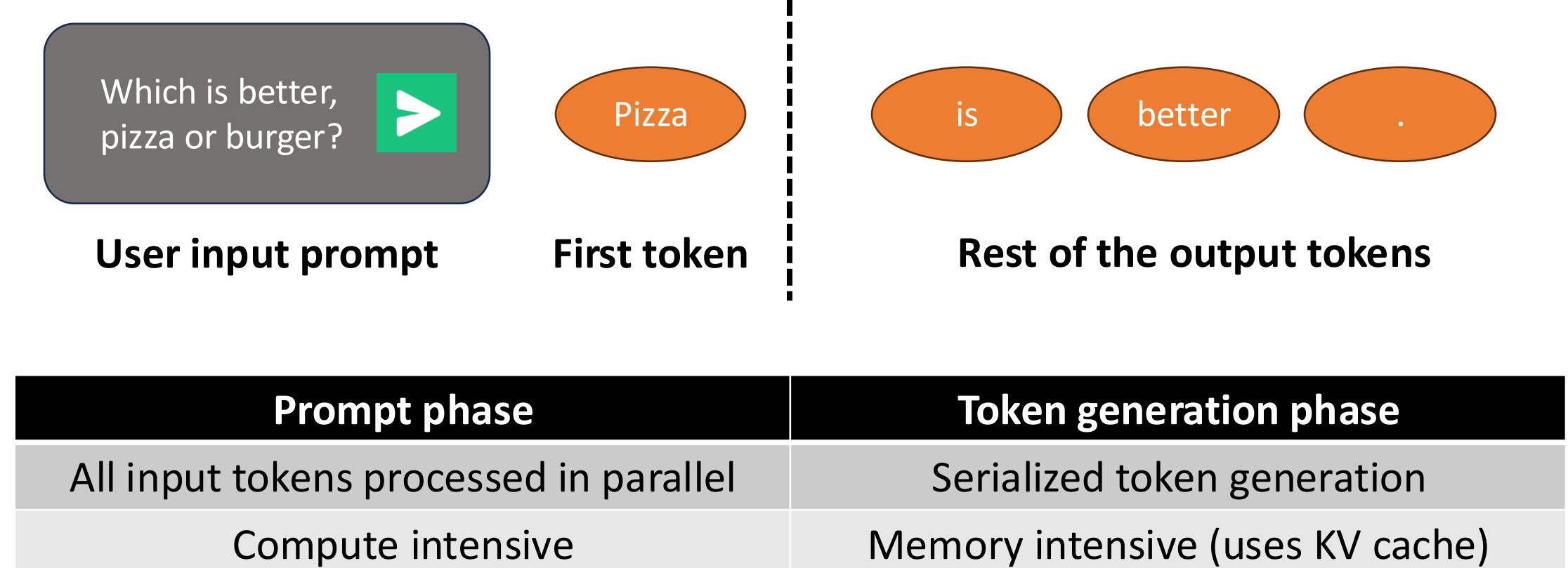
The AI boom is fueling an insatiable appetite for electricity, which is creating risks to the grid and the transition to cleaner energy sources

- The world adds a new datacenter every 3 days
- Datacenter electricity usage set to double by 2026
- 2-to-6 year construction delays due to power supply shortages



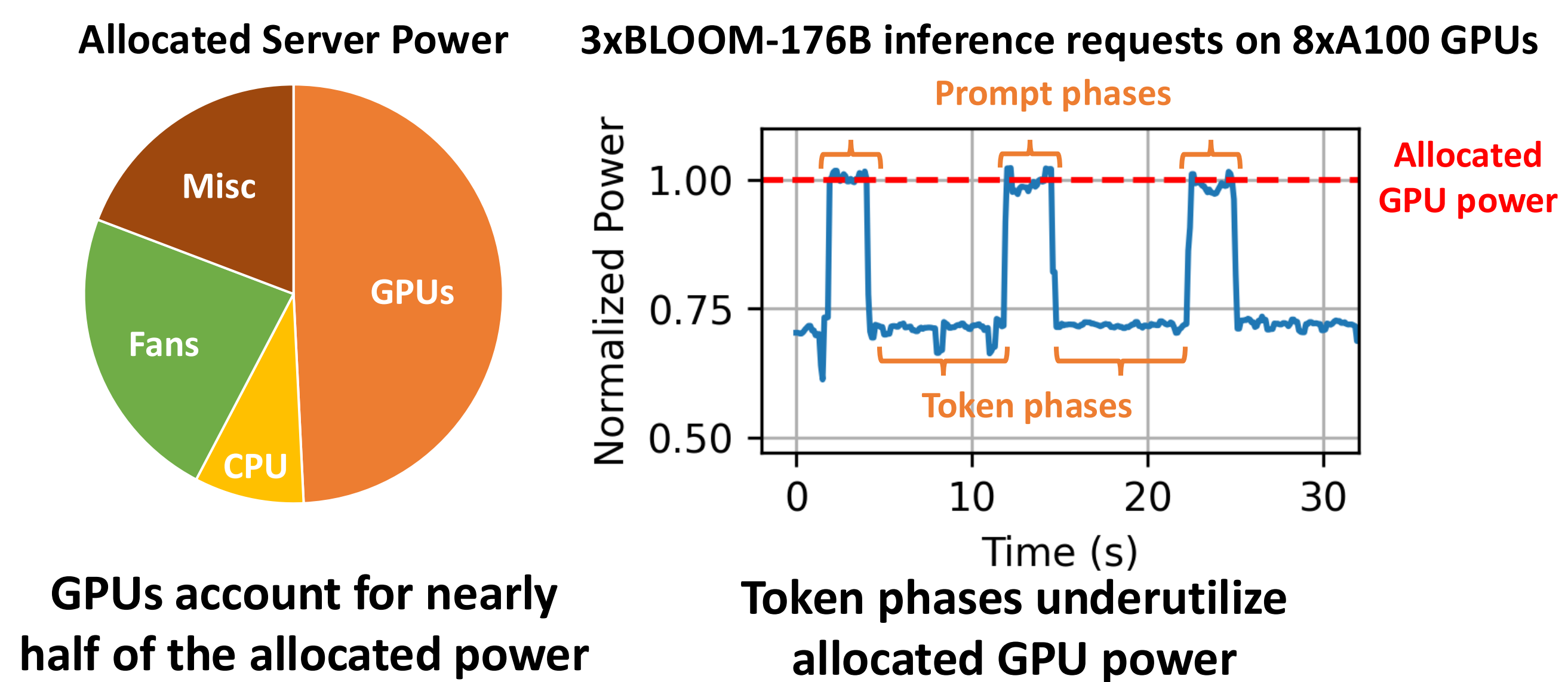
Power supply cannot keep up with the explosion in demand for LLMs

Background: LLM inference has two distinct phases

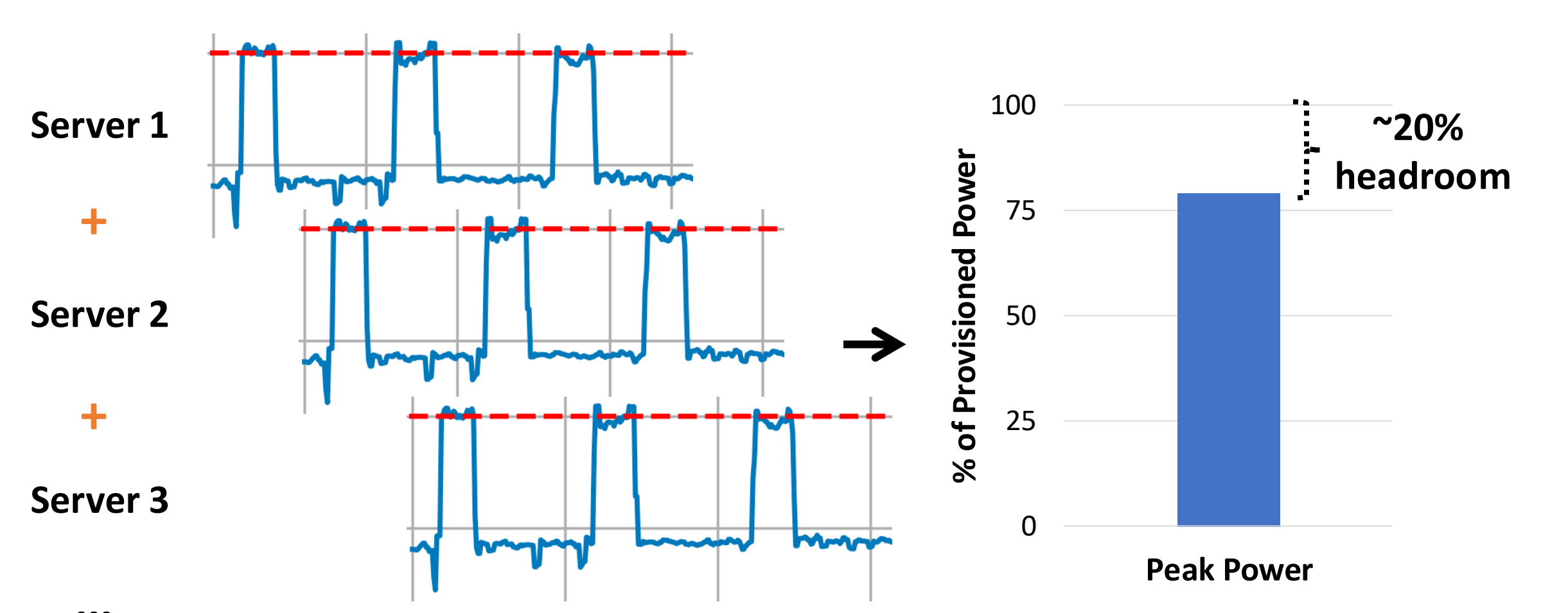


Phases are fundamental to transformer-based generative LLM inference

Prompt phases are power intensive, token phases not

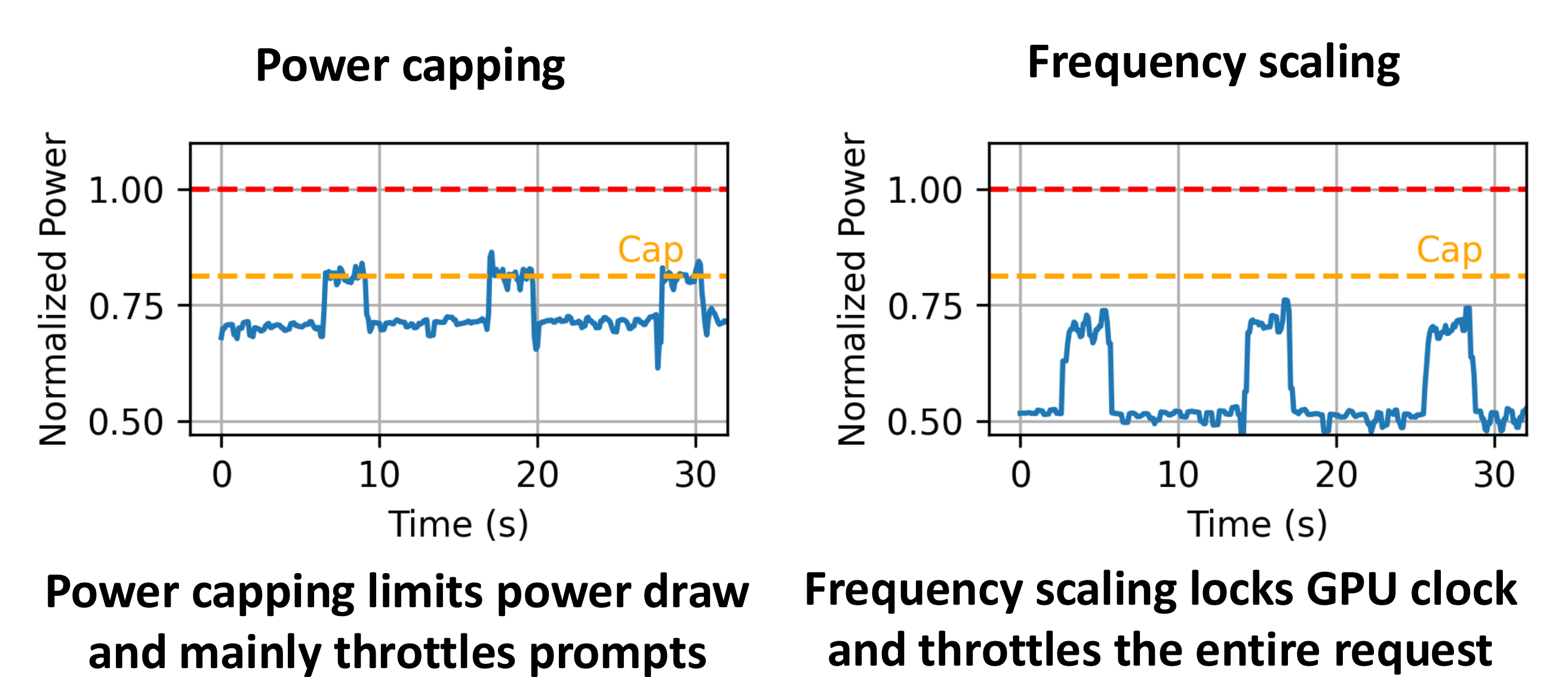


Production LLM inference clusters underutilize power

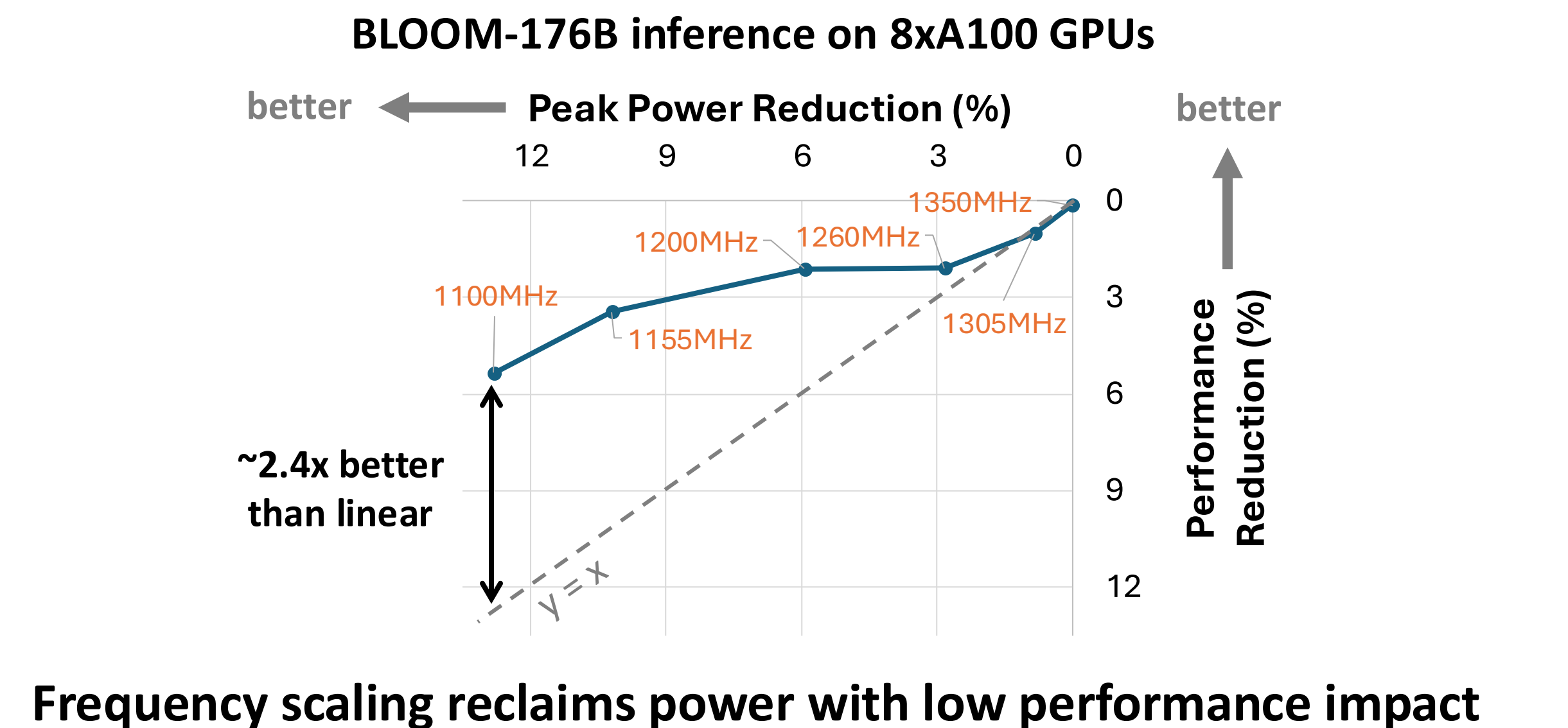


Prompt and token phases are statistically multiplexed across the cluster

GPU power throttling knobs have different trade-offs

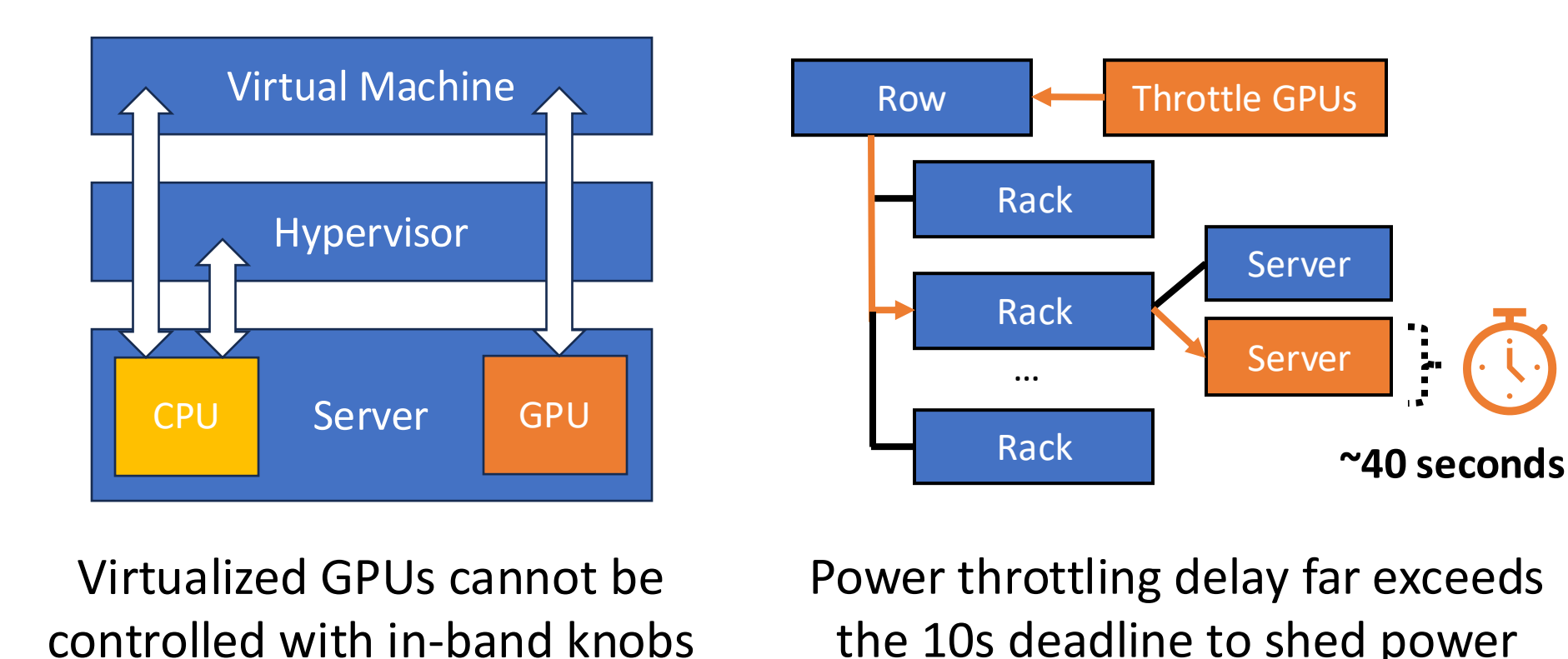


GPU power throttling is effective for LLM inference



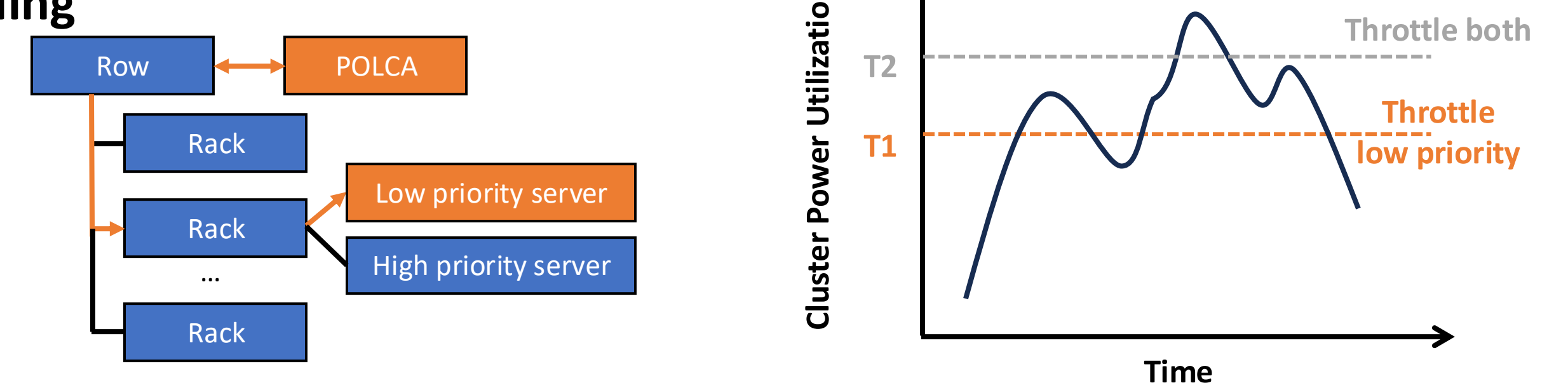
POLCA manages GPU power to safely deploy ~30% more servers in existing and upcoming LLM inference clouds

### Challenge 1: Slow out-of-band GPU management



### Approach: Configurable and proactive priority-aware throttling

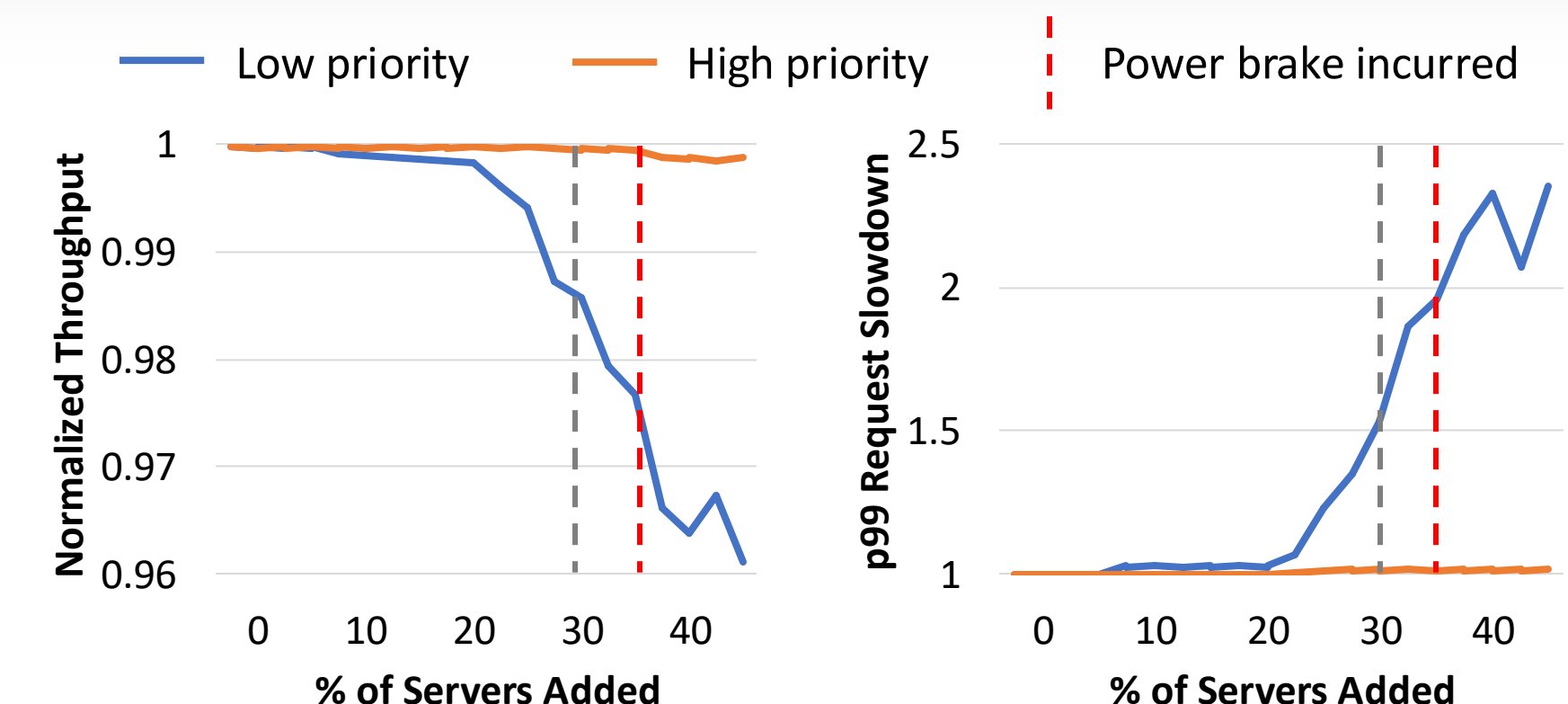
- POLCA Inputs
1. Historical power traces
  2. Workload priorities
  3. Row-level power draw



### Result: Deploy 30% more servers with < 1.5% p99 impact on high-priority workloads

Production power usage patterns with open-source LLMs

Workload	Prompt size	Output size	Fraction
Summarize	2k-8k	256-512	25%
Search	512-2k	1k-2k	25%
Chat	2k-4k	128-2k	50%



### Challenge 2: Diverse and evolving LLMs



In the paper: in-depth training and inference characterization, design implications for LLM clusters, etc.

Pratyush Patel, Esha Choukse, Chaojie Zhang, Íñigo Goiri, Brijesh Warriar, Nithish Mahalingam, Ricardo Bianchini

