

Comparative Analysis of GPU Cloud Providers for AI: Performance, Stability, and Cost-Efficiency

Jihoon Choi, Jehong Bae, Hyeokjin Doo, Young-Woo Kwon
Intelligent Software Systems Lab, Kyungpook National University

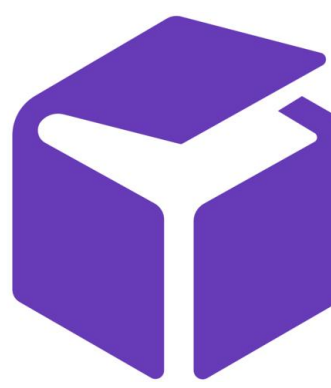
Introduction & Objectives

- The GPU Supply Crisis: The explosive growth of LLM training and serving has led to a global shortage of H100/A100 GPUs, creating a bottleneck for academic and cost-sensitive research.
- The DePIN Alternative: Decentralized Physical Infrastructure Networks (DePINs) utilize idle consumer-grade GPUs (e.g., RTX 4090), offering a potential cost reduction of up to 80% compared to AWS/GCP.
- The Challenge: Unlike centralized clouds, DePINs suffer from hardware heterogeneity, unverified stability, and "noisy neighbor" effects, necessitating a rigorous empirical evaluation.

Vast.ai:
GPU Marketplace
(Direct Rental)



RunPod:
GPU Marketplace/
Dedicated Server
(Direct Rental)



Akash Network:
Decentralized Protocol
(On-chain)



SaladCloud:
Distributed Cloud
(Containerized)



Methodology

To ensure fair comparison across heterogeneous decentralized networks, we developed a portable, automated Python benchmarking suite. As illustrated in Algorithm 1, the framework executes a sequential validation pipeline:

- System Integrity:** Validates GPU availability (CUDA), VRAM capacity, and PCIe bandwidth (H2D/D2H) to detect hardware bottlenecks or PCIe x1 lane limitations.
- Storage & Network:** Measures sequential/random Disk I/O latency and raw network throughput using standard HTTP libraries and HuggingFace Hub downloads.
- AI Inference (vLLM):** Deploys the vLLM engine to measure Tokens Per Second (TPS) and Model Load Time under production-like serving workloads (Algorithm 2).
- Scoring Mechanism:** Raw metrics are normalized against data-center baselines (e.g., 1000 MB/s Disk I/O, 80 TFLOPS) to calculate a weighted performance score (0-100).

Algorithm 1 Unified GPU Cloud Benchmarking Workflow

```
Require: Target Cloud Instance  $I$ 
Ensure: Benchmark Results  $R$ , Score  $S$ 
1:  $INSTALL\_DEPENDENCIES(vLLM, Torch, HF\_Hub)$ 
2:  $H_{info} \leftarrow MEASURE\_SYSTEM\_INFO(Pcie\_Bandwidth, VRAM)$ 
3:  $M_{disk} \leftarrow BENCHMARK\_DISK(SeqWrite, RandRead, Concurrency)$ 
4:  $M_{net} \leftarrow MEASURE\_NETWORK(RawSpeed, ModelDownloadTime)$ 
5: if GPU is Available then
6:    $M_{infer} \leftarrow RUN\_INFERENCE(Model, vLLM)$  ▷ Measure TPS, Load Time
7:    $M_{train} \leftarrow SIMULATE\_TRAINING(GEMM\_Loop)$  ▷ Est. TFLOPS
8: end if
9:  $S \leftarrow CALCULATE\_WEIGHTED\_SCORE(H_{info}, M_{disk}, M_{net}, M_{infer}, M_{train})$ 
10: return  $R \leftarrow \{H_{info}, M_{disk}, M_{net}, M_{infer}, M_{train}\}, S$ 
```

Algorithm 2 Inference Measurement and Scoring Logic

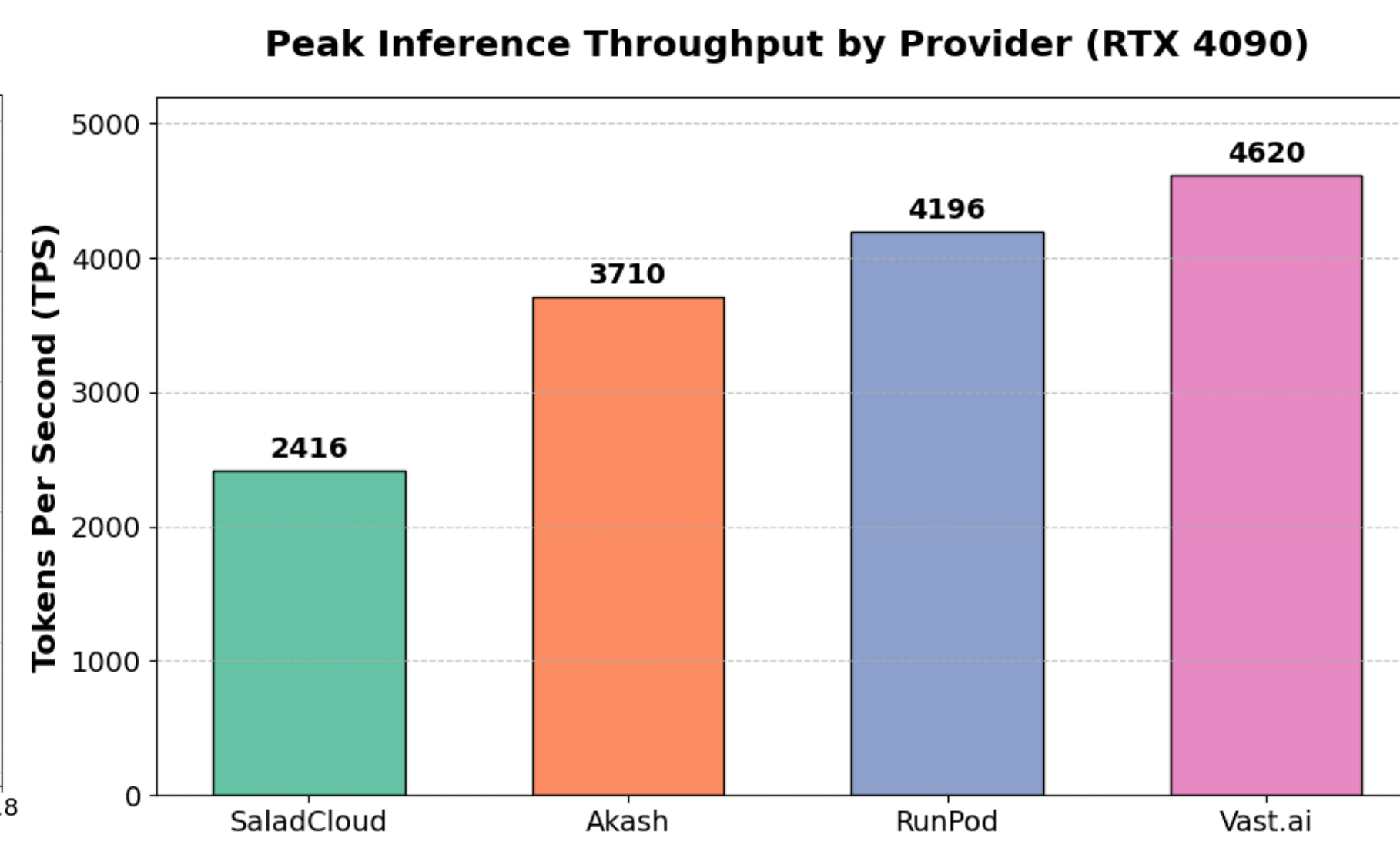
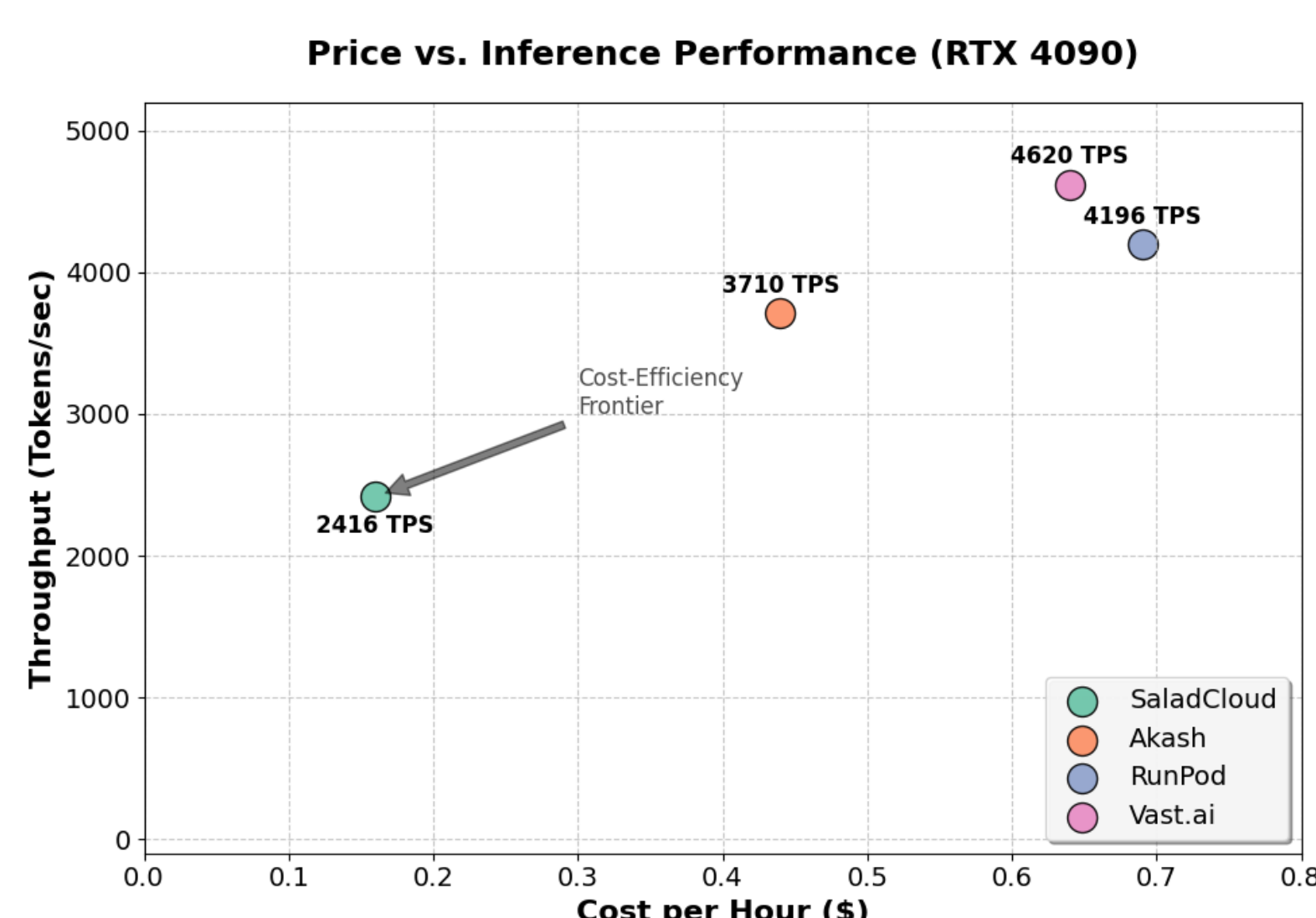
```
1: procedure  $RUN\_INFERENCE(Model)$ 
2:    $t_{start} \leftarrow CURRENT\_TIME$ 
3:    $E \leftarrow LOAD\_ENGINE(Model, vLLM)$ 
4:    $t_{load} \leftarrow CURRENT\_TIME - t_{start}$ 
5:    $O \leftarrow E.GENERATE(Prompts, Params)$ 
6:    $N_{tokens} \leftarrow COUNT\_TOKENS(O)$ 
7:    $TPS \leftarrow N_{tokens} / Duration$ 
8:   return  $TPS, t_{load}$ 
9: end procedure ▷ Score calculation with weighted normalization
10: procedure  $CALCULATE\_SCORE(M)$ 
11:    $S_{train} \leftarrow \min(\frac{M_{train}}{80}, 1.2) \times 50 + \dots$ 
12:    $S_{infer} \leftarrow \min(\frac{M_{infer}}{1000}, 1.2) \times 50 + \dots$ 
13:    $S_{io} \leftarrow \min(\frac{M_{disk}}{1000}, 1.2) \times 40 + \dots$ 
14:   return  $\{S_{train}, S_{infer}, S_{io}\}$ 
15: end procedure
```

Quantitative Performance Summary

Quantitative Performance Summary

Metric	Vast.ai	RunPod	Akash Network	SaladCloud
GPU Model	RTX 4090	RTX 4090	RTX 4090	RTX 4090 / 3090
Cost (\$/hr)	\$0.26 ~ \$0.64	~\$0.69	~\$0.44	\$0.16
Peak Inference (TPS)	4,620	4,196	3,710	2,416
Average PCIe Bandwidth	H2D : 18.99 GB/s D2H : 18.62 GB/s	H2D : 22.21 GB/s D2H : 19.14 GB/s	H2D : 16.30 GB/s D2H : 15.47 GB/s	H2D : 24.30 GB/s D2H : 24.42 GB/s
AVG Seq write (MB/s)	725.262	410.694	302.932	644.745
AVG Small write (MB/s)	163.543	29.917	22.524	83.215
AVG Random read (MB/s)	6021.712	72.303	5292.212	10421.02
AVG Small read (MB/s)	838.352	65.6	1060.499	3498.623

* H2D: Host to Device, D2H: Device to Host



Detailed Provider Analysis

A. RunPod & Vast.ai: The Performance Leaders

- RunPod: Demonstrated exceptional stability and performance with the RTX 4090, achieving 4,196 TPS. Unlike previous trials with other models. However, since runpod utilizes external network drive for workspace, it shows significantly slower disk I/O speed.
- Vast.ai: Achieved the highest theoretical peak (4,620 TPS) on premium instances. However, lower-tier instances (\$0.26/hr) frequently failed due to "Noisy Neighbor" effects (VRAM contention), requiring careful filtering.

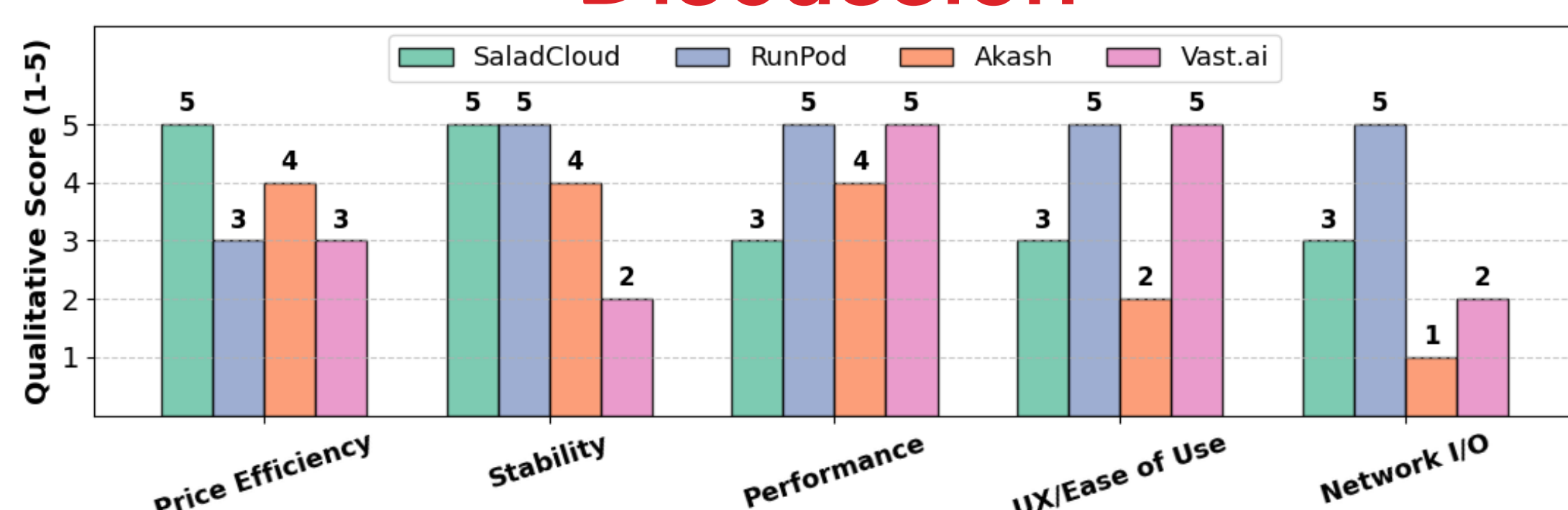
B. Akash Network: The Decentralized Middle Ground

- Observation: Reliable performance (~3,710 TPS) sits comfortably between the high-end and budget options.
- Limitations: Suffers from significant Disk I/O bottlenecks (Seq Write as low as ~30 MB/s in some nodes) and a complex setup process (YAML), reducing rapid deployment capabilities.

C. SaladCloud: The Cost-Efficiency Champion

- Observation: Unbeatable price (\$0.16/hr) with consistent, albeit lower, throughput (TPS ~2,200 - 2,800).
- Trade-off:
 - Provisioning: Long startup times (>10 mins).
 - Access: No SSH access (Web terminal only), limiting advanced debugging.
 - Hardware: Limited to consumer-grade GPUs (No H100/A100).

Discussion



Operational Insights:

- 1. Stability vs. Performance:** While Vast.ai offers the highest peak speeds, RunPod provided a more consistent "out-of-the-box" experience for the RTX 4090 workload without the severe noisy neighbor issues seen in Vast.ai's budget tier.
- 2. IO Variability:** Both RunPod and Vast.ai showed variability in Disk I/O depending on the specific host machine, suggesting that applications requiring heavy data loading should verify storage specs (NVMe vs. Network/SATA) pre-deployment.

Conclusion & Recommendations

The choice of provider depends strictly on the deployment phase

1. For Production (Latency-Critical):

- Recommendation: RunPod or Vast.ai (Premium Tier)
- Reason: Delivers maximum throughput (>4,000 TPS) and full SSH control. RunPod offers a slightly more polished UX, while Vast.ai offers potential cost savings if curated carefully.

2. For Development & Batch Processing:

- Recommendation: SaladCloud
- Reason: extremely low cost makes it perfect for long-running batch jobs where startup latency and SSH are not critical.

3. For Decentralization & Censorship Resistance:

- Recommendation: Akash Network
- Reason: Solid performance with true decentralization, suitable for users who prioritize infrastructure ownership over raw I/O speed.