

# H100 Rental Prices Compared: \$1.49-\$6.98/hr Across 15+ Cloud Providers (2026)

10/24/2025 • 25 min read

[nvidia h100](#)

[gpu rental](#)

[cloud gpu pricing](#)

[ai hardware](#)

[gpu-as-a-service](#)

[aws h100](#)

[machine learning cost](#)



## Executive Summary

By late 2025, renting a NVIDIA H100 80 GB GPU in the cloud has become dramatically cheaper and more competitive than at its 2023 launch. Major cloud providers now offer on-demand H100 rental for **single-GPU rates in the low- to mid-single-digit dollars per hour**, while smaller specialist GPU clouds and marketplaces often charge **well under \$3.00/GPU-hr**. For example, AWS and GCP on-demand H100 pricing stands around \$3–\$4/GPU-hr ([go4hosting.in](#))<sup>[1]</sup> ([www.datacenterdynamics.com](#)), whereas boutique services like Lambda Labs, RunPod, Vast.ai, and Cudo Compute offer rates as low as \$1.49–\$2.99<sup>[2]</sup> ([gpucompare.com](#))<sup>[3]</sup> ([gpucompare.com](#)). These changes follow aggressive price cuts in 2025 (notably AWS cut H100 by ~44% in June 2025<sup>[1]</sup> ([www.datacenterdynamics.com](#))) and **intense competition**.

**Key findings:** The cheapest current H100 rental rates (per GPU-hour) by vendor include:

- **AWS EC2 (P5 instances):** about \$3.90 (per GPU) ([go4hosting.in](#)) on-demand (after mid-'25 price cuts<sup>[1]</sup> ([www.datacenterdynamics.com](#))).
- **Google Cloud (A3-high):** about \$3.00 ([go4hosting.in](#)) on-demand.
- **Microsoft Azure (NC H100 v5):** roughly \$6.98<sup>[4]</sup> ([cloudprice.net](#)) on-demand (in East US).
- **Oracle Cloud (bare-metal 8×H100):** \$10.00 (\$80/hr for 8 GPUs)<sup>[5]</sup> ([www.thundercompute.com](#)).
- **Lambda Labs:** \$2.99<sup>[6]</sup> ([www.thundercompute.com](#)). **CoreWeave:** \$6.16<sup>[7]</sup> ([www.thundercompute.com](#)). **Paperspace:** \$5.95<sup>[8]</sup> ([www.thundercompute.com](#)). **RunPod (community):** \$1.99<sup>[9]</sup> ([www.thundercompute.com](#)). **Vast.ai:** \$1.87 (marketplace)<sup>[10]</sup> ([www.thundercompute.com](#)). **HPC-AI and TensorDock:** \$1.99–\$2.25<sup>[11]</sup> ([gpucompare.com](#)). **Cudo Compute:** \$1.80<sup>[3]</sup> ([gpucompare.com](#)). **VAST.ai** again offers \$1.49 hikes in July 2025<sup>[2]</sup> ([gpucompare.com](#)). Detailed vendor-by-vendor rates are summarized in Tables 1 and 2 below.

These prices vary by region and commitment. Spot/preemptible instances undercut on-demand further (~\$2–2.50/GPU-hr for AWS/GCP ([go4hosting.in](#))). Early adopters paid much more: Thunder Compute reported AWS H100 at ~\$7.57/GPU-hr in Sept 2025<sup>[12]</sup> ([www.thundercompute.com](#)) (pre-cut), Google at ~\$11.06<sup>[13]</sup> ([www.thundercompute.com](#)). The market has since moved dramatically downward (current AWS ~\$3.9 ([go4hosting.in](#)), Google \$3.0 ([go4hosting.in](#))). NVIDIA's own data center systems and device costs are vastly higher (an 8×H100 DGX system is ~\$300,000<sup>[14]</sup> ([www.cudocompute.com](#))), underscoring the appeal of on-demand rental for many users.

Despite the drop, H100 rent remains a premium compared to older GPUs (A100s are now sub-\$1/GPU-hr open-market<sup>[15]</sup> ([www.thundercompute.com](#))). The sustained price reductions reflect oversupply and competition, as well as long-term commitments and spot markets putting downward pressure<sup>[15]</sup> ([www.thundercompute.com](#)) ([go4hosting.in](#)). We examine the historical context, current pricing data, multi-provider comparisons (with extensive citations), market dynamics, and projected trends below.

## Introduction and Background

The **NVIDIA H100 Tensor Core GPU** (Hopper architecture, 80 GB HBM3) is currently NVIDIA's flagship AI accelerator, designed for **large language model training**, HPC, and cloud AI workloads. Released in 2022–2023, the H100 offered **unprecedented performance** (e.g. multi-Exaflop on DGX systems) and swiftly became the gold standard for cutting-edge AI research ([go4hosting.in](#)). However, its *sticker price* at retail (~\$25–40K/GPU ([cyfuture.cloud](#))) makes ownership difficult for many, driving demand for rental solutions.

In parallel, the **GPU-accelerated cloud computing market** has **exploded** in recent years. The global “GPU-as-a-Service” market was only \$3.34 billion in 2023 but is projected to reach \$33.9 billion by 2032<sup>[15]</sup> ([www.thundercompute.com](#)). This growth is driven by **AI/ML needs**, making high-end GPUs like the H100 in continuously tight supply. Industry analysts note that “GPU prices remain sky-high in 2025” and that renting is often the only viable

option for smaller players ([cyfuture.cloud](#)). As adoption grows, vendors have been incentivized to lower prices: Thunder Compute reports H100 rental rates **falling from \$8/hr to \$2.85–3.50/hr** in 2025 (<sup>[15]</sup> [www.thundercompute.com](#)).

**Rental models:** Cloud providers typically charge by the hour (or minute) for GPU instances. Prices depend on *instance type*, region, and commitment. On-demand (no-commitment) is highest; spot/preemptible instances can be 60–90% cheaper (<sup>[16]</sup> [cloud.google.com](#)) ([go4hosting.in](#)); and 1–3 year commitments (Reserved/Savings Plans) offer up to ~45–50% further discounts (<sup>[17]</sup> [aws.amazon.com](#)) ([go4hosting.in](#)). In practice, many users mix and match: e.g. burst on-demand or spot when needed, while using reserved instances for steady-state workloads. The market spans major hyperscalers (AWS, Google, Azure, etc.), specialist AI/cloud vendors (Lambda Labs, CoreWeave, RunPod, etc.), and even AI-focused marketplaces ([Vast.ai](#), etc.).

This report compiles the latest data (as of Nov 2025) on **hourly H100 rental costs across all major and emerging vendors**, analyzes factors behind the prices, and explores implications. We draw on published pricing sources, market analyses, and vendor documentation, with all claims supported by citations. Key takeaways include pricing comparisons (see Tables 1–2), market trend analysis, case examples, and outlook for how H100 pricing may evolve.

## Historical Context and Market Trends

### H100 Release and Early Pricing

When NVIDIA announced the H100 in mid-2022 and it entered general availability in 2023 (<sup>[18]</sup> [aws.amazon.com](#)), it was the fastest GPU for AI training. At launch, public on-demand rental pricing was extremely high. For instance, in mid-2023 AWS's P5 instances (8×H100) were often listed above \$60/hr (~\$7.50/GPU-hr) (<sup>[19]</sup> [www.thundercompute.com](#)) and Google's A3 instances around \$88/hr (~\$11/GPU-hr) (<sup>[13]</sup> [www.thundercompute.com](#)). These figures match NVIDIA's historical pattern of high launch prices, as seen with prior A100/A6000 tiers. Academia & industry narratives at the time noted that while H100 power justified its cost-performance, *cost remained a major barrier* for widespread use.

However, by 2025 the supply of H100s increased substantially (due to added production, alternate chips, and enterprise-focused deals). Cloud providers have also advanced specialized hardware (e.g. AWS Trainium, Google TPUs) which competes with H100. To maintain competitiveness, many began cutting prices. A notable turning point occurred in June 2025: AWS announced a **~44% price reduction** on P5 instances (H100) across regions (<sup>[17]</sup> [aws.amazon.com](#)) (<sup>[1]</sup> [www.datacenterdynamics.com](#)). This brought AWS H100 GPU rental to roughly half its former rate. Similar cuts followed at other clouds. A DataCenterDynamics report confirms: *"On-demand price ... 44% for the H100 instance"* (<sup>[1]</sup> [www.datacenterdynamics.com](#)).

Beyond hyperscalers, many smaller GPU-cloud vendors instituted competition-driven price wars. Throughout 2025, daily tracking sites (e.g. GPUCompare) documented multiple providers slashing H100 rates (see Section on "Daily Pricing Updates" below). By late 2025, aggregated analyses expect H100 rental to settle around **\$2–4/GPU-hr** in most markets (<sup>[15]</sup> [www.thundercompute.com](#)) ([go4hosting.in](#)). For perspective, ThunderCompute's September 2025 snapshot still showed Azure H100 at \$6.98 and GCP at \$11.06 (<sup>[12]</sup> [www.thundercompute.com](#)), but as of fall 2025 current rates are much lower ([go4hosting.in](#)).

key market trend: **price volatility and downward drift**. Indicators include GPUCompare's mid-2025 reports of providers continuously lowering H100 prices by up to 20–25% in single moves (<sup>[3]</sup> [gpucompare.com](#)) (<sup>[2]</sup> [gpucompare.com](#)). Reported insights note *"continued price reductions across board"* (<sup>[20]</sup> [gpucompare.com](#)) and H100 reaching *"new lows"* (<sup>[21]</sup> [gpucompare.com](#)). The H100 price volatility surpasses that of previous-gen GPUs, reflecting fierce competition and excess capacity. Regionally, even within the same cloud provider, pricing differs; for example, the U.S. West Coast regions are consistently ~10–30% above East US rates (<sup>[22]</sup> [www.linkedin.com](#)). Users can often reduce costs by choosing slower regions or leveraging preemptible instances.

## Market Projections

Looking forward, experts anticipate this downward trend to continue, albeit with fluctuations due to demand surges. The ThunderCompute analysis projects global GPU rental market growth from \$3.34B (2023) to \$33.91B by 2032 (<sup>[15]</sup> [www.thundercompute.com](http://www.thundercompute.com)), which should force pricing down over time. The imminent arrival of next-gen GPUs (e.g. NVIDIA H200) may further depress H100 prices as customers upgrade. On the other hand, sustained AI hype could keep demand high, partially offsetting price drops. For now, the evidence points to H100 rentals easing toward a competitive equilibrium of roughly \$2–3/GPU-hr for on-demand usage (<sup>[15]</sup> [www.thundercompute.com](http://www.thundercompute.com)) ([go4hosting.in](http://go4hosting.in)).

## Pricing by Provider

We compared on-demand H100 prices across major cloud and specialized vendors. **Table 1** lists large cloud providers (AWS, Azure, Google, Oracle) based on latest available pricing, normalized per GPU-hour. **Table 2** covers key GPU-cloud specialists and marketplaces. All values are USD per single H100 GPU-hour, using official on-demand rates (no spot). Citations indicate authoritative sources (pricing pages, third-party monitors).

Provider	Instance / SKU	GPUs/Instance	Price (\$/GPU-hr)	Source
AWS ( <a href="http://amazon.com">amazon.com</a> )	p5.48xlarge (8×H100)	8	~\$3.93	( <a href="http://go4hosting.in">go4hosting.in</a> ) (us-east-1)
Google Cloud	A3-highgpu-1g (1×H100)	1	~\$3.00	( <a href="http://go4hosting.in">go4hosting.in</a> ) (us-central)
Azure	Standard_NC40ads_H100_v5 (1×H100)	1	~\$6.98	( <sup>[4]</sup> <a href="http://cloudprice.net">cloudprice.net</a> ) (East US)
Oracle Cloud	BM.GPU.H100.8 (8×H100, bare-metal)	8	\$10.00	( <sup>[5]</sup> <a href="http://www.thundercompute.com">www.thundercompute.com</a> ) (8×\$80/hr)

*Note:* AWS prices reflect the June 2025 reduction (<sup>[1]</sup> [www.datacenterdynamics.com](http://www.datacenterdynamics.com)), so ~\$3.93/GPU-hr (P5.48x in US East)†. Google rates are per on-demand H100 in standard VM. Azure's rate is for a single-GPU NC H100 v5 in East US (<sup>[4]</sup> [cloudprice.net](http://cloudprice.net)). Oracle's bare-metal H100 offering lists \$80/hr for 8 GPUs (<sup>[5]</sup> [www.thundercompute.com](http://www.thundercompute.com)) (i.e. \$10.00/GPU-hr). All values are approximate and may vary by region and time.

Vendor / Service	Instance / Config	GPUs	Price (\$/GPU-hr)	Source
Lambda Labs	Lambda Cloud 8×H100 (SXM80GB)	8	\$2.99	( <sup>[6]</sup> <a href="http://www.thundercompute.com">www.thundercompute.com</a> )
CoreWeave	8×H100 (HGX, IB)	8	\$6.16	( <sup>[7]</sup> <a href="http://www.thundercompute.com">www.thundercompute.com</a> )
Paperspace	Dedicated H100 (80GB)	1	\$5.95	( <sup>[8]</sup> <a href="http://www.thundercompute.com">www.thundercompute.com</a> )
RunPod (community)	Single H100 (80GB PCIe)	1	\$1.99	( <sup>[9]</sup> <a href="http://www.thundercompute.com">www.thundercompute.com</a> )
Vast.ai	H100 (various hosts, marketplace low)	1	~\$1.87	( <sup>[9]</sup> <a href="http://www.thundercompute.com">www.thundercompute.com</a> )
HPC-AI	H100 (SXM)	1	\$1.99	( <sup>[11]</sup> <a href="http://gpucompare.com">gpucompare.com</a> )
TensorDock	H100 SXM5 (dedicated)	1	\$2.25	( <sup>[11]</sup> <a href="http://gpucompare.com">gpucompare.com</a> )
Cudo Compute	H100 (cluster node)	1	\$1.80	( <sup>[3]</sup> <a href="http://gpucompare.com">gpucompare.com</a> )
NeevCloud	H100 (special GPU instance)	1	\$1.79	( <sup>[2]</sup> <a href="http://gpucompare.com">gpucompare.com</a> )
AltCloud etc.	H100 (various offerings)	1	\$1.99–2.79	( <sup>[3]</sup> <a href="http://gpucompare.com">gpucompare.com</a> ) ( <sup>[2]</sup> <a href="http://gpucompare.com">gpucompare.com</a> )

*Notes:* Specialist GPU-cloud vendors typically sell per-GPU instances or clusters. Lambda Labs and CoreWeave give 8×GPU nodes (requiring dividing prices). RunPod's "community cloud" offers an H100 for \$1.99/GPU-hr (<sup>[9]</sup> [www.thundercompute.com](http://www.thundercompute.com)); its "secure cloud" variant is \$2.39 (<sup>[23]</sup> [www.thundercompute.com](http://www.thundercompute.com)). Vast.ai's marketplace finds H100 as cheap as \$1.87 (<sup>[9]</sup> [www.thundercompute.com](http://www.thundercompute.com)). In June–July 2025, competitors aggressively cut rates: e.g.

RunPod to \$2.49–\$2.79 (<sup>[3]</sup> [gpucompare.com](#)), Cudo to \$1.80 (<sup>[3]</sup> [gpucompare.com](#)), Nebius to \$2.00 (<sup>[3]</sup> [gpucompare.com](#)), NeevCloud to \$1.79 (<sup>[2]</sup> [gpucompare.com](#)). The VAST.ai July 2025 sale delivered \$1.49 (<sup>[2]</sup> [gpucompare.com](#)), the current market low.

These figures make clear that **many non-hyperscaler platforms are now cheaper** per GPU-hour than the big three clouds. For instance, Lambda's \$2.99 (<sup>[6]</sup> [www.thundercompute.com](#)) and RunPod's \$1.99 (<sup>[9]</sup> [www.thundercompute.com](#)) compare favorably to **half** of AWS's last MSRP. (AWS's own bulk-reservation prices can approach \$1.90–\$2.10/GPU-hr ([go4hosting.in](#)), but that requires 1–3 year commitments.) Overall, hourly on-demand H100 costs span roughly \$1.50 to \$7.00 across providers, with most falling around \$2–4 by November 2025.

## Detailed Analysis by Vendor

### AWS (Amazon EC2)

AWS offers H100 under its **P5 instance family**. The flagship is `p5.48xlarge` (192 vCPU, 8×H100 80GB). Prior to mid-2025, AWS's on-demand price for `p5.48xlarge` in USD regions was about \$60.54/hour total (<sup>[24]</sup> [www.thundercompute.com](#)) (i.e. \$7.57/GPU-hr) – one of the highest H100 rental rates. However, on June 5, 2025 AWS announced an *“up to 45% price reduction”* for GPU instances, specifically 44% off H100 (P5) on-demand (<sup>[1]</sup> [www.datacenterdynamics.com](#)). Post-cut, the `p5.48xlarge` on-demand went to roughly \$33–\$34/hr (~\$4.1/GPU-hr). According to third-party trackers, current on-demand P5 pricing is about **\$3.90 per GPU-hour** in key regions ([go4hosting.in](#)).

Thus, **AWS now charges roughly \$3.9–4.2/GPU-hr on-demand (US East)** with variations by region and licensing (Windows vs Linux). For example, CloudPrice (Oct 2025) shows `p5.48xlarge` at \$55.04/hr in `us-east-1` (<sup>[25]</sup> [cloudprice.net](#)) (= \$6.88/GPU), possibly due to region or pre-cut data. In general, AWS H100 is no longer the cost leader but remains the most flexible. Note AWS also offers 1–3 year *Savings Plans* which can bring the effective rate below \$2.00 (as low as \$1.90/GPU-hr) ([go4hosting.in](#)) for committed workloads.

### AWS Spot and Reservation Pricing

AWS spot (interruptible) pricing further reduces costs: spot H100 P5 spot rates can drop to roughly **\$2.50/GPU-hr** ([go4hosting.in](#)), depending on region and time of day. Committed Contracts (EC2 Instance Savings Plans) yield ~25-30% more off top of on-demand, so e.g. \$3.90 could become \$1.90–2.10 ([go4hosting.in](#)). This makes AWS competitive for long-term projects, but flexibility is reduced.

### Microsoft Azure

Azure's H100 offering is in the **NC H100 v5 VM family** (NC40ads\_H100\_v5 = 1×H100, 40 vCPU) and the two-GPU **NC80adis\_H100\_v5**. Latest pricing data (Oct 2025) shows **\$6.98/GPU-hr** for a single-H100 NC40 in “East US” (<sup>[4]</sup> [cloudprice.net](#)). Prices vary by region: the lowest is \$6.98 (East US, East US 2, West US 2/3) and can exceed \$9–\$10 in more expensive zones (<sup>[26]</sup> [cloudprice.net](#)). This is notably higher than AWS/GCP. Azure's on-demand cost per GPU remains high because Microsoft has not announced steep cuts for H100 as generically as AWS did.

Box-by-box, Azure's ND and NC families have historically been priced above AWS, reflecting differing market positioning. Users report ~\$7/GPU-hr as typical (<sup>[12]</sup> [www.thundercompute.com](#)). As with AWS, Azure offers 1–3 year Reservations for up to 40–50% discount, but Microsoft's commitment terms are usually around 30% off (depending on region). Azure also has spot (low-priority) VMs for GPUs, though public data on fallback H100 spot rates is scarce as of Nov 2025.

## Google Cloud Platform (GCP)

Google's H100 GPUs come on **A3 VM instances**. The *A3 High GPU* (a3-highgpu-1g) pairs 1×H100 (80GB) with 12 vCPUs, while *A3 Mega* is 4×H100. According to Go4Hosting, the A3 High (1 GPU) on-demand rate is about **\$3.00/GPU-hr** ([go4hosting.in](https://go4hosting.in)). NVIDIA's Spot pricing page confirms preemptible A3-High at **\$2.25/GPU-hr** ([go4hosting.in](https://go4hosting.in)) ([go4hosting.in](https://go4hosting.in)), and A3-Mega at \$2.379/GPU-hr spot ([go4hosting.in](https://go4hosting.in)). Notably, full-price on-demand A3-High is only slightly higher at \$3.00 ([go4hosting.in](https://go4hosting.in)) – Google's sustained-use discounts also lower this further. GCP's on-demand is thus currently among the lowest, slightly undercutting AWS in comparable regions ([go4hosting.in](https://go4hosting.in)).

Historically, GCP was the *most* expensive (with \$11.06/GPU (<sup>[13]</sup> [www.thundercompute.com](https://www.thundercompute.com)) for A3 in Mar 2025), but has rapidly cut prices or rolled out better deals. In practice, Google's flexible billing and sustained-use model can drop costs significantly for continuous use. Preemptibles (analogous to spot) are widely used for training jobs; at \$2.25/GPU-hr they offer huge savings, albeit with occasional interruptions ([go4hosting.in](https://go4hosting.in)).

## Oracle Cloud

Oracle offers H100 GPUs via **bare-metal instances** (the BM.GPU or BMH series). A standard **BM.GPU.H100.8** configuration packs 8×H100 GPUs into one machine. Oracle's posted rate is **\$80.00/hr** for the 8×H100 node (<sup>[5]</sup> [www.thundercompute.com](https://www.thundercompute.com)) (US West), i.e. **\$10.00/GPU-hr**. There is no smaller (1- or 4-GPU) H100 SKU in Oracle's catalog – it is all-or-nothing. This makes Oracle relatively expensive per GPU, though it orders hardware directly with NVIDIA so still offers consistent availability. Oracle notes its bare-metal GPUs waive virtualization overhead, which can be an advantage for HPC workloads, but at a \$10+ price the H100 is a premium in their lineup.

Other large clouds (currently **no official H100 on-demand** from Alibaba Cloud or AWS competitor **Alibaba** has new GPU instances announced, but no public H100 figures as of late 2025). Similarly, **IBM Cloud** started offering H100 in late 2024 but pricing details are not widely published, so we omit them here.

## Lambda, CoreWeave, Paperspace

Lambda Labs, a popular AI cloud, rents H100 GPUs by the unit or in 8×|8 clusters. Lambda's 8×H100 "Lambda Cloud" instances list at **\$2.99/GPU-hr** (<sup>[6]</sup> [www.thundercompute.com](https://www.thundercompute.com)) for SXM (NVL3) usage. (Dividing \$23.92/hr billed for the 8-GPU instance). Individual GPU rentals are available at similar per-GPU rates (i.e. \$2.99 for a node counted per GPU). This is **far below AWS on-demand**.

CoreWeave also offers 8×H100 HGX nodes; their pricing is around \$49.24/hr per node (<sup>[27]</sup> [www.thundercompute.com](https://www.thundercompute.com)), which normalizes to **\$6.16/GPU-hr**. CoreWeave thus is pricier than Lambda, reflecting its niche HPC focus and InfiniBand support.

Paperspace's "dedicated" H100 VM (1×80GB) is **\$5.95/GPU-hr** (<sup>[8]</sup> [www.thundercompute.com](https://www.thundercompute.com)). This is a dedicated single GPU instance, competitive with Azure. Paperspace also provides discounts for multi-GPU bookings, but on-demand single price stays near \$6.

## RunPod, Vast.ai, TensorDock, HPC-AI, Cudo, etc.

Several newer GPU cloud platforms aggressively cut prices. RunPod offers both "community" (unpaired, multi-tenant) and "secure" (dedicated node) H100 rentals. The cheapest RunPod PCIe SKU is **\$1.99/GPU-hr** (<sup>[9]</sup> [www.thundercompute.com](https://www.thundercompute.com)) (community pool price). Its secure H100 is \$2.39, and we note it can vary by machine type.

GPUCompare tracked RunPod lowering H100 PCIe to \$2.49 (<sup>[3]</sup> [gpucompare.com](#)) (with NVL up to \$2.79). So RunPod's market rate is now \$1.99–\$2.79 depending on SKU.

**Vast.ai** is a GPU marketplace where independent hosts list idle GPUs. Vast's lowest H100 listing has been as low as **\$1.87/GPU-hr** (<sup>[9]</sup> [www.thundercompute.com](#)) (used NVIDIA 80GB PCIe GPUs). In July 2025, VAST.ai ran promotions cutting SXM H100 to \$1.49 (<sup>[2]</sup> [gpucompare.com](#)). Vast.ai prices bounce with supply; it often presents the *floor* for what's possible on the open market.

**HPC-AI** (a smaller cloud) advertises H100 at \$1.99 (<sup>[11]</sup> [gpucompare.com](#)). **TensorDock** (Europe/US cloud) similarly offers H100 at \$2.25 (<sup>[11]</sup> [gpucompare.com](#)) (for clients like FloydHub). **Cudo Compute** (a crypto-mining-backed cloud) lists H100 at \$2.45 in many clusters, and was reported to cut to **\$1.80/GPU-hr** in June 2025 (<sup>[3]</sup> [gpucompare.com](#)). Other entrants like **Hyperstack** and **Nebius** have H100 pricing around \$2.00–2.49 after cuts (<sup>[3]</sup> [gpucompare.com](#)).

In summary, *specialist clouds* can be **3–5× cheaper** than the biggest clouds (₹) per-GPU. If one tolerates spot/interruptions or large commitments, H100 hour can dip under \$2 on these platforms (<sup>[3]</sup> [gpucompare.com](#)) ([go4hosting.in](#)). Even among them, prices vary by service level: e.g., RunPod's secure node vs community pools, or Lambda's guaranteed versus spot.

## Spot and Term Discounts

Apart from on-demand, **spot/preemptible** instances further cut costs. GCP's spot H100 is listed at \$2.25 (A3-High) ([go4hosting.in](#)); AWS spot often runs near \$2.50 ([go4hosting.in](#)). Vendor analytics suggest hybrid strategies: use spot to stay "below \$2.5/hr" ([go4hosting.in](#)). Long-term commitments (1–3 year reservations or savings plans) offer deep coupons: AWS users can get effective H100 costs as low as \$1.90–\$2.10 ([go4hosting.in](#)) per GPU-hour (as noted by Go4Hosting), mirroring Intel's previous Frontier writeoff analysis (<sup>[28]</sup> [www.cudocompute.com](#)). GCP sustained-use discounts similarly cut e.g. 20% off after a certain usage threshold.

## Data Analysis and Evidence

Across the data collected from provider sites and third-party trackers, several patterns emerge:

- **Normalized Comparisons:** When comparing per-GPU costs, larger instances distort price-per-GPU. Thus, most analyses (e.g. Thunder Compute) *normalize* to \$/GPU-hr. Table 1 and 2 use this metric.
- **Temporal Trending:** GPUCompare's daily updates document **month-to-month declines**. For example, in June–July 2025 multiple providers slashed H100 pricing by ~20–25% in short order (<sup>[3]</sup> [gpucompare.com](#)) (<sup>[2]</sup> [gpucompare.com](#)). By September, ThunderCompute notes general H100 availability around \$2.85–3.50 (<sup>[15]</sup> [www.thundercompute.com](#)). The cloudfoundation blog reports AWS/Azure now at \$3.9/\$3.0 ([go4hosting.in](#)), down from \$7.57/\$11.06 (<sup>[12]</sup> [www.thundercompute.com](#)). This chronological perspective (Jun 2025: cuts (<sup>[3]</sup> [gpucompare.com](#)) → Sep 2025: ~\$3–4 (<sup>[15]</sup> [www.thundercompute.com](#)) → Nov 2025: similar) shows sustained downward pressure.
- **Regional Variance:** Anecdotal tracking indicates AWS West Coast regions ~15% higher than East (<sup>[22]</sup> [www.linkedin.com](#)). These differences arise from data center costs and local supply. Users can often save by choosing an east/midwest zone. Azure also has up to ~30% spread between cheapest and priciest region (<sup>[29]</sup> [cloudprice.net](#)).
- **Cost vs. Purchase:** For context, the **cost to buy** an H100 is cited at \$25–40K ([cyfuture.cloud](#)). Renting at \$3/GPU-hr full-time (24/7) costs ~\$2,160/month, which would "pay off" the GPU in ~1–2 years. But owning incurs capital, power, facility costs (estimated \$1000–2000/month for electricity alone (<sup>[30]</sup> [www.cudocompute.com](#))). The rental model avoids these and adds maintenance support, which explains why enterprises often prefer cloud rental for short-to-medium-term workloads.
- **Cost-Performance:** Analysts emphasize cost-per-computation, not just sticker price (<sup>[31]</sup> [www.nextplatform.com](#)). For instance, AWS's price/performance improved over prior generations despite higher absolute cost (a theme in The NextPlatform's GPU pricing articles (<sup>[31]</sup>

[www.nextplatform.com](http://www.nextplatform.com)). Our focus here, however, is raw rental price; detailed cost-efficiency (FLOPS per dollar) is beyond this report's scope.

- **Citations and Method:** Table values are drawn from official pricing pages or trusted compilations: e.g. AWS/Azure public calculators, Thunder Compute blog (Sep 2025 data) (<sup>[32]</sup> [www.thundercompute.com](http://www.thundercompute.com)), Go4Hosting analysis ([go4hosting.in](http://go4hosting.in)) ([go4hosting.in](http://go4hosting.in)), GPUCompare daily logs (<sup>[33]</sup> [gpucompare.com](http://gpucompare.com)) (<sup>[3]</sup> [gpucompare.com](http://gpucompare.com)), and vendor docs. Where providers only list multi-GPU rates, we divided by GPU count (noting that licensed per-instance often prices by node). We avoided any hidden fees or long-term commitments in these figures.

## Case Studies and Examples

**Startup GPU procurement:** A CloudCombinator analysis of GenAI startups highlights the complexities of acquiring H100 capacity on AWS (<sup>[34]</sup> [cloudcombinator.ai](http://cloudcombinator.ai)) (<sup>[35]</sup> [cloudcombinator.ai](http://cloudcombinator.ai)). It notes that “popular GPU types (H100, P6-B200) may be out of stock” in busy regions (<sup>[35]</sup> [cloudcombinator.ai](http://cloudcombinator.ai)), forcing planners to use alternatives like GPU Capacity Reservations or AWS UltraClusters for guaranteed access (<sup>[36]</sup> [cloudcombinator.ai](http://cloudcombinator.ai)). This underscores that price alone isn't enough; availability is also a factor. The blog recommends strategies combining spot buys, regional flexibility, and multi-cloud to manage costs—echoing our observation of price variation.

**Cloud vs. On-Premises TCO:** A Cudo Compute blog breaks down the total cost of owning an 8-H100 server. They estimate each H100 costs ~\$30,971 part-price (<sup>[28]</sup> [www.cudocompute.com](http://www.cudocompute.com)), so an 8-card system ~\$247,766 plus CPU (\$25K) and extras (<sup>[37]</sup> [www.cudocompute.com](http://www.cudocompute.com)). They note at least \$300,000 total including power/cooling (<sup>[14]</sup> [www.cudocompute.com](http://www.cudocompute.com)). Renting 8 H100s on Cudo at \$19.60/hour (their quoted multi-GPU rate) costs \$39,292/month, making up-front roughly \$300K in ~7–8 months of 24/7 rental (<sup>[38]</sup> [www.cudocompute.com](http://www.cudocompute.com)). This simple analysis (exclusive of ops costs) suggests renting is more cost-effective for shorter projects. For many, annual or monthly GPU rental provides budget flexibility and eliminates maintenance overhead.

**Regulatory procurement:** Government and academic HPC centers often publish utilization cases. While specific H100 rental deals are proprietary, general budget reports (e.g. Argonne/Lawrence projects) note that hardware budgets hammered by rising chip prices make cloud an attractive option. For example, the Frontier supercomputer's on-site GPU parts were amortized to ridiculously low hourly prices (Argonne essentially paid \$0/GPU-hr after write-off (<sup>[39]</sup> [www.nextplatform.com](http://www.nextplatform.com))). This contrasts with cloud's retail rates and hints why high-demand research might qualify for subsidized purchases instead.

## Implications and Future Directions

### Economic and Adoption Implications

The steep fall in H100 rental costs has several implications. For enterprises and researchers, AI training becomes more accessible: fine-tuning large models or running vision pipelines on H100 is now far cheaper than a year ago. This democratization fuels innovation but also squeezes margins for small GPU vendors. Market analysts observe “**GPU rental market exploded**” demand but concurrently pressure to “drive down costs” (<sup>[15]</sup> [www.thundercompute.com](http://www.thundercompute.com)). In effect, the anticipated GPU shortage of 2023–24 eased somewhat, as vendors responded by lowering prices and increasing supply chains (purchase from multiple chip fabs, interconnect tech improvements, etc).

For cloud providers, these price cuts reflect economies of scale and competitive necessity. AWS's 45% cut was partly driven by customer pushback over scarcity pricing. Google and Microsoft risk losing GPU-hungry customers to smaller clouds, which offer more attractive hourly rates (albeit usually with less enterprise support). Conversely, providers expect to retain workloads by offering discounts for committed usage (1–3 years) – a reason AWS still advertises ~\$1.90/GPU-hr

savings plans ([go4hosting.in](https://go4hosting.in)). Potentially, if H100 prices drop further, providers may shift emphasis to newer chips (like H200) as the premium offering.

## Technological Currency and Inventory

Another aspect is inventory management. Some smaller clouds (e.g. Lambda Labs, RunPod) offering H100 at \$2/GPU-hr risk hardware resale losses if H100 market prices fall below purchase cost. Indeed, Intel's reports show per-GPU cost halving over generations (<sup>[31]</sup> [www.nextplatform.com](https://www.nextplatform.com)). If H100 commodity price falls, companies may offload old inventory cheap. Conversely, manufacturer reactions (e.g. Nvidia's own pricing guidance) could support higher prices via limiting supply of discounts. But as of Nov 2025, H100 chips are no longer "limited edition" – plenty of refurbished or gamer-stock boards have appeared.

## Future GPUs (H200 and beyond)

NVIDIA H100's life cycle will impact its rental price trajectory. Introduced in 2022–23, the next-generation **H200 (Blackwell)** is expected ~2026. With new chips on horizon, H100 may see further price declines as datacenters upgrade. Indeed, some early evidence shows clouds already marketing H200 instances at premium, which effectively hopes to supplant H100 demand. We would expect by mid-2026 that H100 rent might fall into sub-\$2 range universally, and older GPUs (A100/A6000) become nearly free (<\$1). The same logic drove GPUCompare's late-2024 commentary on H100 pricing stability despite volatility (<sup>[20]</sup> [gpucompare.com](https://gpucompare.com)); eventually, as H200 arrives, H100 will enter "previous-generation" pricing logic fast.

## Broader Impact

Falling rent may hasten AI adoption in SMBs and academia. For startups and labs previously deterred by cost, H100 compute is reaching commodity price points. This could contribute to an AI boom or democratize research. However, some caution is warranted: extremely low prices often come with caveats (limited slots, older GPUs sold as H100, or multi-tenancy on slower networking). Additionally, as H100 rent flattens, differences between providers may hinge on value-adds (support, specialized interconnects, compliance). It's also possible that in later 2026, if GPU supply outstrips demand, providers will deliberately throttle lower-end rates to protect profit margins – e.g. by promoting multi-year commitments or preemptible offerings over on-demand.

## Conclusion

In summary, the **H100 GPU rental market by November 2025 is both mature and competitive**. Effective per-GPU-hour rates have fallen to roughly \$2–4 on large clouds and often under \$3 on smaller clouds (with some as low as \$1.50–2.00) ([go4hosting.in](https://go4hosting.in)) (<sup>[2]</sup> [gpucompare.com](https://gpucompare.com)). These figures are backed by multiple sources: cloud pricing pages, market analyst reports, and GPU-compare trackers. We have tabulated the latest on-demand rates (Tables 1–2) to give a definitive snapshot.

Key drivers have been broad market trends – surging supply, price wars, and alternative chips – all documented here with citations (<sup>[1]</sup> [www.datacenterdynamics.com](https://www.datacenterdynamics.com)) (<sup>[3]</sup> [gpucompare.com](https://gpucompare.com)) (<sup>[15]</sup> [www.thundercompute.com](https://www.thundercompute.com)) ([go4hosting.in](https://go4hosting.in)). Users planning ML workloads can expect to pay a few dollars per GPU-hour for H100 compute in late 2025. Those able to leverage spot instances or reserved plans will pay even less (often ~\$2/GPU-hr or below) ([go4hosting.in](https://go4hosting.in)) ([go4hosting.in](https://go4hosting.in)).

Overall, H100 rental has transitioned from a boutique premium to a commoditized resource. It remains a "high-end" GPU, but thanks to intense market forces its price approach established norms. For budgeting purposes: **plan on \$3–6 per**





## IntuitionLabs - Industry Leadership & Services

**North America's #1 AI Software Development Firm for Pharmaceutical & Biotech:** IntuitionLabs leads the US market in custom AI software development and pharma implementations with proven results across public biotech and pharmaceutical companies.

**Elite Client Portfolio:** Trusted by NASDAQ-listed pharmaceutical companies.

**Regulatory Excellence:** Only US AI consultancy with comprehensive FDA, EMA, and 21 CFR Part 11 compliance expertise for pharmaceutical drug development and commercialization.

**Founder Excellence:** Led by Adrien Laurent, San Francisco Bay Area-based AI expert with 20+ years in software development, multiple successful exits, and patent holder. Recognized as one of the top AI experts in the USA.

**Custom AI Software Development:** Build tailored pharmaceutical AI applications, custom CRMs, chatbots, and ERP systems with advanced analytics and regulatory compliance capabilities.

**Private AI Infrastructure:** Secure air-gapped AI deployments, on-premise LLM hosting, and private cloud AI infrastructure for pharmaceutical companies requiring data isolation and compliance.

**Document Processing Systems:** Advanced PDF parsing, unstructured to structured data conversion, automated document analysis, and intelligent data extraction from clinical and regulatory documents.

**Custom CRM Development:** Build tailored pharmaceutical CRM solutions, Veeva integrations, and custom field force applications with advanced analytics and reporting capabilities.

**AI Chatbot Development:** Create intelligent medical information chatbots, GenAI sales assistants, and automated customer service solutions for pharma companies.

**Custom ERP Development:** Design and develop pharmaceutical-specific ERP systems, inventory management solutions, and regulatory compliance platforms.

**Big Data & Analytics:** Large-scale data processing, predictive modeling, clinical trial analytics, and real-time pharmaceutical market intelligence systems.

**Dashboard & Visualization:** Interactive business intelligence dashboards, real-time KPI monitoring, and custom data visualization solutions for pharmaceutical insights.

**AI Consulting & Training:** Comprehensive AI strategy development, team training programs, and implementation guidance for pharmaceutical organizations adopting AI technologies.

Contact founder Adrien Laurent and team at <https://intuitionlabs.ai/contact> for a consultation.

---

## DISCLAIMER

The information contained in this document is provided for educational and informational purposes only. We make no representations or warranties of any kind, express or implied, about the completeness, accuracy, reliability, suitability, or availability of the information contained herein.

Any reliance you place on such information is strictly at your own risk. In no event will IntuitionLabs.ai or its representatives be liable for any loss or damage including without limitation, indirect or consequential loss or damage, or any loss or damage whatsoever arising from the use of information presented in this document.

This document may contain content generated with the assistance of artificial intelligence technologies. AI-generated content may contain errors, omissions, or inaccuracies. Readers are advised to independently verify any critical information before acting upon it.

All product names, logos, brands, trademarks, and registered trademarks mentioned in this document are the property of their respective owners. All company, product, and service names used in this document are for identification purposes only. Use of these names, logos, trademarks, and brands does not imply endorsement by the respective trademark holders.

IntuitionLabs.ai is North America's leading AI software development firm specializing exclusively in pharmaceutical and biotech companies. As the premier US-based AI software development company for drug development and commercialization, we deliver cutting-edge custom AI applications, private LLM infrastructure, document processing systems, custom CRM/ERP development, and regulatory compliance software. Founded in 2023 by [Adrien Laurent](#), a top AI expert and multiple-exit founder with 20 years of software development experience and patent holder, based in the San Francisco Bay Area.

This document does not constitute professional or legal advice. For specific guidance related to your business needs, please consult with appropriate qualified professionals.

© 2025 IntuitionLabs.ai. All rights reserved.