

GPU As A Service - Market Share Analysis, Industry Trends & Statistics, Growth Forecasts (2026 - 2031)

Market Report | 2026-01-16 | 137 pages | Mordor Intelligence

AVAILABLE LICENSES:

- Single User License \$4750.00
- Team License (1-7 Users) \$5250.00
- Site License \$6500.00
- Corporate License \$8750.00

Report description:

GPU As A Service Market Analysis

GPU as a Service market size in 2026 is estimated at USD 7.36 billion, growing from 2025 value of USD 5.70 billion with 2031 projections showing USD 26.43 billion, growing at 29.12% CAGR over 2026-2031. The GPU as a Service market draws momentum from the collision of generative-AI workloads, cloud-gaming adoption, and companywide digital-transformation projects that require elastic, high-density compute capacity. Pay-per-use models continue to shift budgets away from on-premises GPU clusters toward cloud subscriptions, while liquid-cooling retrofits enable data-center operators to pack more accelerators per rack and maintain power efficiency. Hyperscalers protect share through global scale, yet specialist "neoclouds" compete aggressively on price and workload-specific performance. Pricing ranges from USD 0.66 per hour for A100 instances to USD 4.00 and above for premium H100 configurations, giving customers flexibility across performance tiers.

Global GPU As A Service Market Trends and Insights

Rising usage of generative-AI and LLM workloads

Demand for transformer-based models drives unprecedented GPU clustering, with single projects consuming thousands of H100 accelerators for training cycles that last weeks. NVIDIA noted that 91% of financial institutions are now in production or evaluation phases for AI use cases. Financial-services firms such as BNY Mellon demonstrated the power of GPU superclusters for real-time fraud analytics nvidia.com. Elastic scaling inherent in the GPU as a Service market allows research teams to match compute

supply with unpredictable training bursts. High-bandwidth memory (HBM) equipped H100 and H200 parts are favored because they maintain throughput for expanding parameter counts. The long tail of startups can now access the same silicon that hyperscalers deploy, leveling the innovation playing field.

Surge in AR/VR and real-time rendering needs

Photorealistic rendering at 90 frames per second strains consumer hardware, motivating developers to stream pixel-perfect frames from remote GPUs. NVIDIA's CloudXR platform layers low-latency codecs onto GPU back-ends to deliver immersive experiences to thin clients. Pixel-streaming specialists such as Arcware offer Unreal-Engine-as-a-Service so that architectural-visualization teams can present interactive models on mobile devices. Manufacturing firms adopt digital-twin workflows that mix physics simulation with real-time visualization, pushing demand for distributed GPUs at the edge. As next-generation headsets arrive, content studios prefer the GPU as a Service market over purchasing bespoke render farms because they avoid capital costs and maintain flexibility.

Cyber-security and data-sovereignty concerns

Shared accelerator pools create fresh attack surfaces, with research highlighting GPU side-channel vectors that bypass traditional hypervisor barriers. Confidential-computing extensions now encrypt memory and isolate workloads so that multi-tenant environments meet bank and government standards. Export-control regimes add compliance complexity because GPUs above certain TOPS thresholds require licensing before cross-border deployment. Sovereign-cloud frameworks push enterprises toward regional nodes, influencing data-center location strategies inside the GPU as a Service market. Providers respond with per-region key-management systems and cryptographically signed GPU-license enforcement.

Other drivers and restraints analyzed in the detailed report include:

Cloud-gaming service expansionPay-per-use pricing models gaining tractionHBM memory and advanced packaging supply constraints

For complete list of drivers and restraints, kindly check the Table Of Contents.

Segment Analysis

Artificial-intelligence use cases represented 46.68% of 2025 revenue, giving this segment the largest slice of the GPU as a Service market. Transformer architectures now exceed 1 trillion parameters, driving multi-cluster demands that only elastic cloud pools can supply. Large-language-model inference spans real-time chatbots, code-generation assistants, and enterprise knowledge retrieval, keeping utilization steady after training cycles complete.

Cloud Gaming and Media Rendering is the fastest-rising application group at a 30.35% CAGR, helping expand the GPU as a Service market size for entertainment workloads through 2031. Providers monetize evening gaming peaks and rent idle daytime capacity to film-render pipelines, elevating asset utilization. Hybrid workloads that simulate autonomous-vehicle environments blend photoreal rendering with physics-based AI, bridging gaming engines and AI frameworks in a single tenancy. As these cross-domain workflows normalize, application boundaries blur and every incremental project funnels additional value into the GPU as a Service market.

Large Enterprises secured 55.54% of 2025 revenue thanks to reserved-capacity contracts and dedicated support teams. Multinational banks, automakers, and pharmaceutical giants lock in multi-year blocks of H100 instances for predictable AI roadmaps. They often negotiate data-center colocation arrangements or direct-to-manufacturer supply guarantees, ensuring

continuity during supply-chain shocks.

Small and Medium Enterprises are growing at a 28.95% CAGR, underscoring the democratization effect that consumption billing brings to the GPU as a Service industry. Serverless offerings remove the need for DevOps headcount, allowing lean teams to integrate vision models or recommendation engines into products within days. Competitive pricing at USD 0.66 per hour for A100s further lowers entry barriers, propelling the overall GPU as a Service market forward as SME participation deepens.

GPU As A Service Market is Segmented by Application (Artificial Intelligence, High-Performance Computing, and More), Enterprise Size (Small and Medium Enterprises, Large Enterprises), End-User Industry (BFSI, Automotive and Mobility, and More), Deployment Model (Public Cloud, Private Cloud, and Hybrid / Multi-Cloud), Service Model (IaaS, PaaS, and More), and by Geography. The Market Forecasts are Provided in Terms of Value (USD).

Geography Analysis

North America contributed 30.88% of global revenue in 2025 on the back of established hyperscaler footprints, vibrant startup ecosystems, and early adoption across banking, retail, and entertainment. Providers retrofit legacy halls with direct-to-chip liquid-cooling to achieve rack densities above 120 kW, enabling tens of thousands of GPUs per facility. Regional export controls shape where the most advanced silicon can be deployed, adding compliance consulting as a value-added service inside the GPU as a Service market.

Asia-Pacific is projected to post a 29.70% CAGR, reflecting government-funded AI clouds and manufacturing digitization. Singapore spends USD 600 per capita on NVIDIA hardware and offers tax incentives for AI infrastructure, positioning itself as a regional compute hub. India's national mission to install 10,000 GPUs partners NVIDIA with domestic telcos for sovereign-cloud builds. Japan and South Korea accelerate procurement of H200 clusters for language-translation and robotics workloads, illustrating diverse catalysts that funnel regional budgets into the GPU as a Service market.

Europe balances growth opportunities with stringent sustainability and data-residency regulations. Providers invest in 100% renewable energy supplies and waste-heat re-use, aligning with EU carbon caps. Demand across automotive, pharma, and public-sector AI applications keeps utilization rising despite policy headwinds. Growth in South America and the Middle East & Africa lags in absolute terms but posts double-digit gains as broadband penetration improves and local AI ecosystems mature. Collectively, emerging regions will expand the addressable user base and further diversify revenue streams for the GPU as a Service market.

List of Companies Covered in this Report:

Amazon Web Services Microsoft Azure NVIDIA DGX Cloud Google Cloud IBM Cloud Oracle Cloud Alibaba Cloud CoreWeave Linode / Akamai Latitude.sh Seeweb Lambda Labs Paperspace (DigitalOcean) Vultr OVHcloud Scaleway RunPod Vast.ai Genesis Cloud Cirrascale

Additional Benefits:

 The market estimate (ME) sheet in Excel format
3 months of analyst support

Table of Contents:

1 INTRODUCTION
1.1 Study Assumptions and Market Definition

Scotts International. EU Vat number: PL 6772247784

tel. 0048 603 394 346 e-mail: support@scotts-international.com
www.scotts-international.com

1.2 Scope of the Study

2 RESEARCH METHODOLOGY

3 EXECUTIVE SUMMARY

4 MARKET LANDSCAPE

4.1 Market Overview

4.2 Market Drivers

4.2.1 Rising usage of generative-AI and LLM workloads

4.2.2 Surge in AR/VR and real-time rendering needs

4.2.3 Cloud-gaming service expansion

4.2.4 Pay-per-use pricing models gaining traction

4.2.5 Liquid-cooling data-center retrofits unlocking GPU density

4.2.6 Multi-cloud GPU-orchestration platforms reducing vendor lock-in

4.3 Market Restraints

4.3.1 Cyber-security and data-sovereignty concerns

4.3.2 Global shortage of AI-skilled DevOps talent

4.3.3 HBM memory and advanced packaging supply constraints

4.3.4 Escalating data-center power tariffs and carbon regulations

4.4 Supply-Chain Analysis

4.5 Regulatory Landscape

4.6 Technological Outlook

4.7 Porter's Five Forces Analysis

4.7.1 Bargaining Power of Suppliers

4.7.2 Bargaining Power of Buyers

4.7.3 Threat of New Entrants

4.7.4 Threat of Substitutes

4.7.5 Intensity of Competitive Rivalry

4.8 Assesment of Macroeconomic Factors on the market

5 MARKET SIZE AND GROWTH FORECASTS (VALUE)

5.1 By Application

5.1.1 Artificial Intelligence

5.1.2 High-Performance Computing

5.1.3 Cloud Gaming and Media Rendering

5.1.4 Other Applications

5.2 By Enterprise Size

5.2.1 Small and Medium Enterprises

5.2.2 Large Enterprises

5.3 By End-user Industry

5.3.1 BFSI

5.3.2 Automotive and Mobility

5.3.3 Healthcare and Life Sciences

5.3.4 IT and Communications

5.3.5 Media and Entertainment

5.3.6 Other Industries

5.4 By Deployment Model

5.4.1 Public Cloud

5.4.2 Private Cloud

5.4.3 Hybrid / Multi-cloud

5.5 By Service Model

5.5.1 IaaS

5.5.2 PaaS

5.5.3 SaaS (GPU-accelerated)

5.6 By Geography

5.6.1 North America

5.6.1.1 United States

5.6.1.2 Canada

5.6.1.3 Mexico

5.6.2 South America

5.6.2.1 Brazil

5.6.2.2 Argentina

5.6.2.3 Rest of South America

5.6.3 Europe

5.6.3.1 United Kingdom

5.6.3.2 Germany

5.6.3.3 France

5.6.3.4 Italy

5.6.3.5 Spain

5.6.3.6 Rest of Europe

5.6.4 Asia-Pacific

5.6.4.1 China

5.6.4.2 Japan

5.6.4.3 South Korea

5.6.4.4 India

5.6.4.5 Australia

5.6.4.6 Rest of Asia-Pacific

5.6.5 Middle East and Africa

5.6.5.1 Middle East

5.6.5.1.1 Saudi Arabia

5.6.5.1.2 United Arab Emirates

5.6.5.1.3 Turkey

5.6.5.1.4 Rest of Middle East

5.6.5.2 Africa

5.6.5.2.1 South Africa

5.6.5.2.2 Nigeria

5.6.5.2.3 Egypt

5.6.5.2.4 Rest of Africa

6 COMPETITIVE LANDSCAPE

6.1 Market Concentration

6.2 Strategic Moves

6.3 Market Share Analysis

Scotts International. EU Vat number: PL 6772247784

tel. 0048 603 394 346 e-mail: support@scotts-international.com

www.scotts-international.com

6.4 Company Profiles (includes Global level Overview, Market level overview, Core Segments, Financials as available, Strategic Information, Market Rank/Share for key companies, Products and Services, Recent Developments)

6.4.1 Amazon Web Services

6.4.2 Microsoft Azure

6.4.3 NVIDIA DGX Cloud

6.4.4 Google Cloud

6.4.5 IBM Cloud

6.4.6 Oracle Cloud

6.4.7 Alibaba Cloud

6.4.8 CoreWeave

6.4.9 Linode / Akamai

6.4.10 Latitude.sh

6.4.11 Seeweb

6.4.12 Lambda Labs

6.4.13 Paperspace (DigitalOcean)

6.4.14 Vultr

6.4.15 OVHcloud

6.4.16 Scaleway

6.4.17 RunPod

6.4.18 Vast.ai

6.4.19 Genesis Cloud

6.4.20 Cirrascale

6.5 Vendor Ranking Analysis

7 MARKET OPPORTUNITIES AND FUTURE OUTLOOK

7.1 White-space and Unmet-need Assessment

Scotts International. EU Vat number: PL 6772247784

tel. 0048 603 394 346 e-mail: support@scotts-international.com

www.scotts-international.com

GPU As A Service - Market Share Analysis, Industry Trends & Statistics, Growth Forecasts (2026 - 2031)

Market Report | 2026-01-16 | 137 pages | Mordor Intelligence

To place an Order with Scotts International:

- Print this form
- Complete the relevant blank fields and sign
- Send as a scanned email to support@scotts-international.com

ORDER FORM:

Select license	License	Price
	Single User License	\$4750.00
	Team License (1-7 Users)	\$5250.00
	Site License	\$6500.00
	Corporate License	\$8750.00
		VAT
		Total

*Please circle the relevant license option. For any questions please contact support@scotts-international.com or 0048 603 394 346.

** VAT will be added at 23% for Polish based companies, individuals and EU based companies who are unable to provide a valid EU Vat Numbers.

Email*	<input type="text"/>	Phone*	<input type="text"/>
First Name*	<input type="text"/>	Last Name*	<input type="text"/>
Job title*	<input type="text"/>		
Company Name*	<input type="text"/>	EU Vat / Tax ID / NIP number*	<input type="text"/>
Address*	<input type="text"/>	City*	<input type="text"/>
Zip Code*	<input type="text"/>	Country*	<input type="text"/>
		Date	<input type="text" value="2026-02-07"/>

Signature

Scotts International. EU Vat number: PL 6772247784

tel. 0048 603 394 346 e-mail: support@scotts-international.com

www.scotts-international.com



Scotts International. EU Vat number: PL 6772247784

tel. 0048 603 394 346 e-mail: support@scotts-international.com

www.scotts-international.com