

# Private AI Solutions: A Guide to Datacenter Providers

By IntuitionLabs.ai • 10/23/2025 • 35 min read

private ai

ai infrastructure

generative ai

data center providers

on-prem ai

hybrid cloud

gpu computing



## Executive Summary

The rapid advent of generative AI has driven an unprecedented surge in demand for specialized data center infrastructure. Global AI spending is estimated at **\$1.5 trillion by 2025**, with **AI-optimized servers** (\$267 B in 2025) and supporting chips (\$209 B) dominating infrastructure investment (<sup>[1]</sup> [www.gartner.com](http://www.gartner.com)). As hyperscalers (e.g. Amazon Web Services, Microsoft Azure, Google Cloud) expand ultra-powerful GPU-based data centers (<sup>[2]</sup> [www.gartner.com](http://www.gartner.com)), many enterprises are simultaneously **shifting back toward private/hybrid solutions**. Data privacy, security and cost predictability concerns have led CIOs to revive on-premises and colocation "private AI" deployments (<sup>[3]</sup> [www.cio.com](http://www.cio.com)) (<sup>[4]</sup> [www.cio.com](http://www.cio.com)). Leading providers from all categories – cloud hyperscalers, hardware vendors, virtualization platforms, and colocation operators – now offer "private AI" offerings: fully isolated, GPU-rich infrastructure hosted on-premises or in dedicated facilities. For example, Equinix advertises its global colocation footprint as "the place where private AI happens," enabling companies to leverage public AI models while keeping proprietary data off the public Internet (<sup>[5]</sup> [www.equinix.com](http://www.equinix.com)) (<sup>[6]</sup> [www.equinix.com](http://www.equinix.com)). Major hardware vendors are partnering with NVIDIA to bundle turnkey AI data center stacks (HPE's *Private Cloud AI*, Cisco's *Nexus HyperFabric AI*, Dell's GPU-optimized servers, etc.) (<sup>[7]</sup> [www.datacenterknowledge.com](http://www.datacenterknowledge.com)) (<sup>[8]</sup> [newsroom.cisco.com](http://newsroom.cisco.com)) (<sup>[9]</sup> [www.reuters.com](http://www.reuters.com)). Enterprise software players like VMware and IBM have introduced private-AI architectures (e.g. *VMware Private AI Foundation*, IBM watsonx on-prem) to run customized LLMs in-house (<sup>[10]</sup> [www.techtarget.com](http://www.techtarget.com)) (<sup>[11]</sup> [www.ibm.com](http://www.ibm.com)). At the same time, security vendors (e.g. Trend Micro) stress private-cloud AI deployments for **compliance and data protection** (<sup>[12]</sup> [newsroom.trendmicro.com](http://newsroom.trendmicro.com)).

This report provides a comprehensive analysis of the **leading data center providers of private-AI solutions** and their offerings. We examine the historical context of enterprise AI computing (moving from on-prem to public cloud and now to hybrid/private), characterize private-AI solutions across multiple provider categories, and present detailed case studies. We compare specific products and platforms (hyperscale cloud services, on-prem appliances, colocation offerings) and analyze industry data on spending and adoption. Finally, we discuss broader implications – from energy consumption of AI data centers (<sup>[13]</sup> [www.allaboutai.com](http://www.allaboutai.com)) to future trends (confidential computing, sovereign clouds) – and conclude with an outlook on where the private-AI market is headed. All assertions are supported by current industry reports, news stories, and expert commentary.

## Introduction and Background

The recent breakthroughs in generative AI (large language models, multimodal vision models, etc.) have dramatically **escalated the need for extreme computational resources**. Training and fine-tuning state-of-the-art AI models can require thousands of high-end GPUs working in parallel, as well as ultra-fast networking and storage. According to Gartner, global spending on AI (infrastructure, software, and services) is projected to reach **\$1.48 trillion in 2025** (<sup>[1]</sup> [www.gartner.com](http://www.gartner.com)), with **AI-optimized servers** (high-density GPU clusters) alone consuming roughly **\$267 billion** in 2025. A Distinguished Gartner analyst notes, "major hyperscalers continue to increase investments in data centers with **AI-optimized hardware and GPUs** to scale their services" (<sup>[2]</sup> [www.gartner.com](http://www.gartner.com)). This has led to major capital commitments: **Microsoft**, for example, allotted **\$80 billion in fiscal 2025** to build data centers designed for model training and AI services (<sup>[14]</sup> [www.reuters.com](http://www.reuters.com)). Likewise, industry consortia (e.g. OpenAI/Oracle/SoftBank's "**Stargate**" project) plan massive new AI training farms.

Initially, much of this compute has been supplied by **cloud hyperscalers** (AWS, Azure, Google Cloud, etc.), which have the economies of scale to quickly deploy new GPU generations (NVIDIA H100, H200, AMD MI300X, etc.) (<sup>[15]</sup> [azure.microsoft.com](http://azure.microsoft.com)) (<sup>[16]</sup> [azure.microsoft.com](http://azure.microsoft.com)). However, enterprises are increasingly wary of hosting **highly sensitive or proprietary data** on shared cloud platforms, even with encryption. As Paula Rooney (CIO magazine) notes, **AI data leak fears are driving CIOs to rethink cloud strategies**: many now plan a hybrid mix,

using private clouds for critical workloads (<sup>[3]</sup> [www.cio.com](http://www.cio.com)) (<sup>[4]</sup> [www.cio.com](http://www.cio.com)). In fact, IDC predicts that by 2025, Global 2000 firms will spend over **40% of their core IT budgets on AI initiatives** (<sup>[17]</sup> [www.equinox.com](http://www.equinox.com)), and a large portion of those initiatives will demand **private, isolated infrastructure** for compliance and control.

The term *Private AI* (or *private cloud AI*) refers to deployments where the entire AI stack – from base models to training data – runs on dedicated infrastructure (company-owned datacenters, on-premises servers, or leased single-tenant facilities) rather than a multi-tenant cloud. This model offers **full isolation** (“nobody else’s tenants or ‘noisy neighbors’ share your GPUs”) (<sup>[18]</sup> [www.nexgencloud.com](http://www.nexgencloud.com)), which helps prevent data leakage and reduces latency variability. Enterprises adopting private AI cite benefits such as **complete data control, predictable costs, and ease of regulatory compliance**. For example, one bank executive explained that while public clouds have the horsepower for many LLMs, the bank prefers to keep critical data in a private environment (using on-prem Dell GPU servers) to avoid any chance of it getting ingested into a third-party model (<sup>[19]</sup> [www.cio.com](http://www.cio.com)).

This resurgence of private/hybrid infrastructure is reminiscent of the earlier “private cloud” wave in the 2010s. However, AI workloads have *greater scale, security, and networking demands* than typical enterprise apps. As IDC analyst Peter Rutten observes, “AI has different system, data, and privacy requirements than existing workloads” (<sup>[20]</sup> [www.datacenterknowledge.com](http://www.datacenterknowledge.com)). Companies thus need not only raw GPU power, but also specialized data transfer (e.g. NVIDIA Quantum InfiniBand), high-throughput storage (GPUDirect Storage), and tightened security (enclave/sealed deployments). The emerging market of “private AI solutions” is therefore quite broad, involving hyperscaler clouds extending on-premises, hardware vendors bundling HPC clusters, virtualization/hybrid-cloud platforms adding AI, and data center companies provisioning GPU racks.

Below, we examine each category of leading providers and detail their private-AI offerings. We cover (a) **hyperscale clouds** (who still support hybrid AI), (b) **IT hardware and system vendors** (HPE, Dell, etc.), (c) **virtualization/middleware providers** (VMware, IBM, etc.), (d) **colocation/datacenter operators** (Equinix, Digital Realty, etc.), and (e) **other solution partners** (security firms, specialized AI hosts). We include specific product names, performance specs, and pricing models where available. Throughout, we cite surveys, case studies, and analyst commentary to ground our discussion in data and expert opinion.

## Hyperscale Cloud Providers and Hybrid AI

**Amazon Web Services (AWS), Microsoft Azure, and Google Cloud** remain the dominant providers of AI infrastructure, but each has introduced hybrid or private-cloud options tailored for AI workloads:

- **AWS:** While AWS excels at public cloud AI (e.g. [SageMaker](#) managed ML service, Bedrock LLM service, etc.), it also offers on-prem solutions. *AWS Outposts* and *AWS Local Zones* bring AWS compute (including GPU instances) closer to customer sites. In 2024 AWS launched **ChatX**, a “private GenAI platform” for enterprise customers in Thailand – essentially a turnkey ChatGPT-like chatbot running in the customer’s own AWS account (<sup>[21]</sup> [www.blognone.com](http://www.blognone.com)). AWS also supports AMD MI300X accelerator VMs (as an alternative to NVIDIA) (<sup>[22]</sup> [www.reuters.com](http://www.reuters.com)). However, AWS has not publicly branded a “private AI” program akin to VMware or HPE. In practice, enterprises can use AWS’s dedicated hardware (e.g. *Nitro Enclaves*, *Graviton3 Pro* servers) and networking to build secure private AI clouds, but solutions are often custom engagements rather than off-the-shelf.

- **Microsoft Azure:** Microsoft has aggressively positioned Azure for AI workloads. At Ignite 2023, Microsoft unveiled the **Azure Maia** and **Azure Cobalt** custom chips, along with updated network fabrics, to supply massive in-cloud AI compute ([16] azure.microsoft.com) ([23] azure.microsoft.com). Azure's strategy also explicitly supports hybrid/edge AI: the **Azure Arc** platform lets customers deploy Azure AI services and Kubernetes clusters on their own servers or other clouds ([24] azure.microsoft.com). For example, Azure Arc now enables running VMware vSphere workloads (including AI/ML) under Azure management ([24] azure.microsoft.com). In practice, an enterprise could deploy Azure Stack Hub (on-prem Azure), connect it via Azure Arc, and run Azure's AI tools locally. Microsoft also sells the **Azure Stack Edge** appliance (GPU-accelerated) for on-site AI inference. On the software side, Azure offers private LLM capabilities such as **Azure OpenAI Service with virtual network isolation**. Microsoft's ongoing investment (\$80B in FY2025 ([14] www.reuters.com)) shows its focus on in-cloud AI; but its hybrid offerings make it a key private-AI provider as well.
- **Google Cloud (GCP):** Google markets *Google Distributed Cloud (GDC)* as a turnkey on-prem/edge solution for AI workloads ([25] cloud.google.com). GDC hardware comes preinstalled in customer sites (data centers or edge locations) with NVIDIA GPUs (now H100-based) and is managed by Google. In Dec 2024 Google launched a "Gen AI Search" packaged solution on GDC: it includes a private-instance LLM (Gemma 2) and connectors to on-prem data, letting enterprises run conversational search locally ([26] cloud.google.com). The GDC appliances support air-gapped operation if needed ([25] cloud.google.com). Google also sells *Anthos*, a Kubernetes-based hybrid platform that can host ML workloads across cloud and on-prem. In short, Google provides both the hardware platform (GDC servers) and software stacks for private AI deployments, though adoption among enterprises is still emerging.
- **Others (Oracle, Alibaba, IBM, etc.):** Oracle Cloud Infrastructure (OCI) offers dedicated regions (OCI Dedicated Region) for enterprises to have Oracle-managed cloud hardware on-prem. While not marketed specifically as "private AI," these can host AI workloads under Oracle's umbrella. Alibaba Cloud provides GPU instances and has been expanding AI infrastructure (e.g. next-gen AI chips), but it is chiefly a Chinese cloud with less presence in global corp. **IBM Cloud** distinguishes itself with **Watsonx** – IBM's AI platform – which can run on-prem via VMware or Red Hat OpenShift. In 2024 IBM announced a partnership with VMware to run Watsonx on private clouds (VCF/OpenShift) ([11] www.ibm.com), explicitly targeting genAI use cases behind the company firewall. IBM also offers AI-ready *Power Systems* servers greenlit for NVIDIA AI Enterprise software.

In summary, while AWS/Azure/GCP lead in raw AI compute capacity, each major cloud player provides hybrid on-prem options (appliances, stack extensions, turnkey deployments) to address private AI use cases. These offerings often integrate the same GPUs and software (NVIDIA AIE, Kubernetes, etc.) found in public clouds, but are delivered as isolated installations.

## Enterprise Hardware and Systems Vendors

Hardware manufacturers and system integrators have rapidly developed **turnkey AI-infrastructure solutions** for enterprises wanting on-prem GPUs:

- **Hewlett Packard Enterprise (HPE):** HPE's flagship **Private Cloud AI** is a fully integrated stack co-developed with NVIDIA ([7] www.datacenterknowledge.com) ([27] www.datacenterknowledge.com). Announced in June 2024, it bundles HPE ProLiant servers, high-speed Ethernet (NVIDIA Spectrum-X), storage (GreenLake file), and NVIDIA AI Enterprise software, all managed via HPE GreenLake (cloud control plane) ([7] www.datacenterknowledge.com). Private Cloud AI comes in four configurations: from small 4-GPU setups (NVIDIA L40S) up to large 8–16 GPU systems using H100/H200 or the new GH200 "Grace Hopper" superchips ([27] www.datacenterknowledge.com). It specifically targets generative AI inference (RAG search) and fine-tuning, with embedded tagging of customer data via a lakehouse ([28] www.datacenterknowledge.com) ([29] www.datacenterknowledge.com). HPE also introduced specialized DL380/DL384 ProLiant servers and Cray XD670 for GPU training. Analysts view HPE Private Cloud AI as a comprehensive on-premises AI platform, leveraging GreenLake's easy consumption model ([29] www.datacenterknowledge.com) ([27] www.datacenterknowledge.com).

- **Dell Technologies:** Dell is introducing its own AI-ready hardware. In May 2025 Dell announced **new AI servers** equipped with NVIDIA's latest Blackwell Ultra GPUs (<sup>[9]</sup> [www.reuters.com](http://www.reuters.com)). These modular systems support up to 192 GPUs (and will support future NVIDIA Vera CPUs/GB200 chips) for "four times faster" training performance (<sup>[9]</sup> [www.reuters.com](http://www.reuters.com)). Dell offers these servers air-cooled or liquid-cooled. Dell also promotes its **APEX** portfolio for hybrid cloud, including APEX Gateway and Cloud Platform (for Azure), which can run on-prem sitting adjacent to AI workloads. Dell's strong presence in high-performance compute earned it a reported **\$5 billion AI server deal with Elon Musk's xAI** (to supply GB200-based clusters) (<sup>[30]</sup> [www.reuters.com](http://www.reuters.com)), highlighting enterprise appetite for dedicated GPU racks.
- **Cisco Systems:** Known mainly for networking, Cisco now sells turnkey AI clusters. The **Cisco Nexus HyperFabric AI Cluster** (announced June 2024) marries Cisco's high-end Ethernet switches (6000 series, 400/800Gb) with NVIDIA GPUs and DPUs in a prevalidated rack design (<sup>[8]</sup> [newsroom.cisco.com](http://newsroom.cisco.com)) (<sup>[31]</sup> [newsroom.cisco.com](http://newsroom.cisco.com)). This solution is cloud-managed via Cisco's Networking Cloud, offering a single pane to design/deploy/monitor an on-prem AI fabric (<sup>[32]</sup> [newsroom.cisco.com](http://newsroom.cisco.com)) (<sup>[31]</sup> [newsroom.cisco.com](http://newsroom.cisco.com)). Key components include NVIDIA Tensor Core GPUs (H200 NVL), BlueField-3 DPUs, NVIDIA AI Enterprise software/NIM microservices, and VAST Data storage (<sup>[31]</sup> [newsroom.cisco.com](http://newsroom.cisco.com)). Cisco's pitch is a simplified "plug-and-play" AI data center: enterprises get a reference design built on Cisco Silicon One, with automated deployment tools and end-to-end visibility. IDC notes this approach leverages Cisco's portfolio strength, positioning them strongly in enterprise AI switching and infrastructure (<sup>[33]</sup> [newsroom.cisco.com](http://newsroom.cisco.com)).
- **VMware (Broadcom):** VMware has framed *Private AI* as a cloud architecture. In August 2023 VMware announced **VMware Private AI** (with NVIDIA) – a stack of VMware Cloud Foundation (vSphere, NSX, vSAN) plus NVIDIA AI Enterprise (including NIM inference services) (<sup>[34]</sup> [www.techtarget.com](http://www.techtarget.com)) (<sup>[10]</sup> [www.techtarget.com](http://www.techtarget.com)). The goal is to let enterprises train LLMs on-prem while preserving VMware's manageability. For example, VMware demonstrated an LLM running on up to 16 virtual GPUs in one VM (<sup>[35]</sup> [www.techtarget.com](http://www.techtarget.com)). Private AI Foundation became generally available in 2024, bundled as part of VCF and including optimizations for bare-metal performance with GPU pass-through (<sup>[10]</sup> [www.techtarget.com](http://www.techtarget.com)). VMware's carbon copy is to eventually have cloud providers host *VMware Private AI* via their managed SDDC offerings. Meanwhile, the Dell, HPE, and Lenovo server partners are expected to certify VCF-based Private AI stacks.
- **IBM:** IBM leverages its enterprise AI stack *watsonx*. In late 2023, IBM announced a collaboration with VMware to enable **watsonx on-prem** (VMware Private AI on OpenShift) (<sup>[11]</sup> [www.ibm.com](http://www.ibm.com)). This means customers can run IBM's AI software (including Watsonx AI Studio and Watsonx data catalog) in dedicated VMware-powered private clouds, with optional IBM Cloud Satellite for management. The emphasis is on *confidential computing* and governance for sensitive AI workloads. IBM also sells on-prem GPUs via its PowerAI platform and the IBM Cloud Private (via Kyndryl) for hybrid AI. In practice, IBM's offerings blur into consulting services: e.g. IBM Consulting has a generative AI CoE with 1,000+ experts to tailor private AI deployments.
- **Other hardware vendors:** Lenovo, Supermicro, and other server makers also offer high-density GPU machines. For example, Supermicro's GPU servers are popular in HPC colo facilities. **NVIDIA MIC** (Modular AI Compute Platform) and **HPE Cray XD** target OEMs building AI farms. However, compared to the above, these vendors typically supply components rather than full solutions. (Dell/Lenovo/HPE OEM many Supermicro/RedHat).

In sum, **hardware and system vendors** have quickly assembled integrated on-prem AI data centers. All point to tight NVIDIA partnerships: "if we go to a customer and don't put NVIDIA in front of them, they will walk away," notes an IDC analyst (<sup>[36]</sup> [www.datacenterknowledge.com](http://www.datacenterknowledge.com)). The result is a hot market for "AI racks" – complete with NVIDIA GPUs, NVLink networking, and accompanying software (NVIDIA AI Enterprise, CUDA libraries, MLOps tools) – marketed as ready-to-run genAI infrastructure. These solutions often include managed service options (via GreenLake, APEX, Cisco/VMware subs), meaning companies can essentially "subscribe" to a private AI data center.

## Virtualization and Hybrid-Cloud Providers

Software platforms that span private/public clouds are also key enablers of private AI:

- **VMware:** Beyond the Private AI Foundation stack (described above), VMware's core **vSphere/vSAN/NSX** technology underpins many private-AI solutions. VMware Private AI is often sold as part of VMware Cloud Foundation on-prem. VMware expects large enterprises will consume the Private AI software on their own hardware. Partners (Dell, HPE, Lenovo) are packaging VCF+NVIDIA stacks. VMware also highlights integration with **NVIDIA AI Enterprise** (for Kubernetes, etc.) and Nebula GPT (their own LLM inference accelerator). In announcements it promised vGPU performance "bare-metal-fast" with vMotion capability (<sup>[37]</sup> [www.techtarget.com](http://www.techtarget.com)). In effect, VMware sells the abstraction layer that lets IT teams treat an on-prem GPU cluster like a private cloud with familiar tools.
- **Red Hat / OpenShift:** Red Hat (owned by IBM) offers Kubernetes-based private cloud (OpenShift) which supports AI workloads. OpenShift AI add-ons and the **Red Hat Open Data Hub** provide Jupyter, ML pipelines, and model serving on-prem. Red Hat's main value is lock-in reduction: you can run the same container orchestration on AWS, Azure, or on-prem (via OpenShift operator) for AI apps. IBM's watsonx on-prem is designed to run on OpenShift.
- **NVIDIA:** NVIDIA themselves have become a platform provider. Their **NVIDIA AI Enterprise** software (containerized frameworks, NIM microservices, Morpheus security framework, etc.) is licensed for private data centers. NVIDIA also operates the **AI LaunchPad** program (with Equinix, HPE, etc.) to get enterprises immediate GPU clusters on demand (<sup>[38]</sup> [www.equinix.com](http://www.equinix.com)). Moreover, NVIDIA has started its own cloud offering (NVIDIA DGX Cloud) that can connect enterprise data centers securely, and **NGC** for model management. In the private domain, any NVIDIA DGX system (A100 or H100) is effectively a building block of a private AI box.
- **Lenovo:** Lenovo partners with NVIDIA as well (DGX-Ready rack solutions) and offers "AI Foundry" servers. It was one of the first Broadcom partners to announce VMware Private AI on its ThinkSystem servers. It also supplies high-end Blade/ThinkEdge nodes to colos. Specific offerings aren't widely publicized, but Lenovo's channel is active in private cloud deals.
- **Other Software:** Some companies provide end-to-end platforms that include private AI. For example, **NexGen Cloud** (a UK-based specialist) markets a "Private AI Cloud" service – essentially dedicated GPU clusters in Tier-3 data centers with full isolation and InfiniBand networking (<sup>[39]</sup> [www.nexgencloud.com](http://www.nexgencloud.com)) (<sup>[40]</sup> [www.nexgencloud.com](http://www.nexgencloud.com)). Startups like **Lambda** and **CoreWeave** (GPU cloud providers) also let customers deploy private instances of GPUs in their facilities (often colocated). Emerging "GPU over IP" solutions (e.g. Juice Labs) are another model: enterprises get virtual private GPU pools via leased hardware at colocation sites.

## Colocation and Data Center Providers

Colocation firms and data center operators have added AI services to attract enterprise GPU deployments. The most prominent:

- **Equinix, Inc.:** Equinix brands itself as a natural home for private AI. In a December 2023 press release, Equinix stated that its *Platform Equinix* ("cloud adjacency, global reach, robust ecosystem") makes it a "preferred location for deploying private AI infrastructure" (<sup>[41]</sup> [www.equinix.com](http://www.equinix.com)). Equinix argues that AI workloads benefit from low-latency interconnection and proximity to cloud on-ramps, preventing sensitive data from leaving the enterprise (<sup>[5]</sup> [www.equinix.com](http://www.equinix.com)). IDC agrees, noting Equinix's network density is well-suited to private AI needs (<sup>[42]</sup> [www.equinix.com](http://www.equinix.com)). Equinix has even dubbed itself "the place where private AI happens;" according to a company executive (<sup>[6]</sup> [www.equinix.com](http://www.equinix.com)).

Equinix has collaborated with NVIDIA via the **NVIDIA AI LaunchPad** (since 2021) to install ready-to-use GPU racks for customers (<sup>[38]</sup> [www.equinix.com](http://www.equinix.com)). It also runs a marketplace of GPU service providers: case studies include Crusoe Energy (mobile gas-powered GPU pods connected via Equinix fabrics), Juice Labs (GPU-over-IP service leveraging Equinix's footprint), and Lambda (GPU cloud) (<sup>[43]</sup> [www.equinix.com](http://www.equinix.com)). Equinix's own IBX data centers have been used in multiple private-AI customer projects: for instance, Continental AG (automotive) trains vision models in Equinix with NVIDIA and IBM Storage (<sup>[44]</sup> [www.equinix.com](http://www.equinix.com)), and Harrison.ai (healthtech) deploys NVIDIA DGX A100 units in Equinix Sydney for medical imaging models (<sup>[45]</sup> [www.equinix.com](http://www.equinix.com)). The Equinix press release cites IDC projections that by 2025 Fortune 2000 firms will shift an increasing share of IT spend to AI (<sup>[17]</sup> [www.equinix.com](http://www.equinix.com)) – a strategic tailwind. In short, Equinix offers the space, power, and

interconnect (direct links to public clouds, partners, users) that enterprises need to build **federated or fully private AI grids**.

- **Digital Realty (PlatformDIGITAL):** Digital Realty (owner of Interxion, Telecity) similarly markets its global platform for private AI/ML. Its *PlatformDIGITAL*® provides secure data center campuses with high-speed interconnect and Software Defined Exchanges for cloud adjacency. While less vocal on AI than Equinix, Digital Realty emphasizes hybrid-cloud flexibility. It recently discussed a “Private AI Exchange” concept for high-speed, secure data sharing around AI workloads. The company participates in projects like NVIDIA-Ready colocation but details are scarcer. However, given its 300+ global data centers, any enterprise could host GPU racks in DR sites and peer to clouds.
- **Other Colos and Carrier Clouds:** Many regional carriers and colo providers now offer AI-specific services. For example, **Kyndryl** (spun out from IBM) manages dedicated AI clusters in partner centers. Companies like **LightEdge** (US) and **Node4** (UK) advertise “GPU colocation” packages. Asia-Pacific players (e.g. NTT Global Data Centers, KDDI) similarly tout high-density AI-ready floors. In China, local data center operators work closely with Baidu/Tencent for private AI networks. In Europe, sovereign or telco clouds (OVHcloud, Orange Flexible Engine, etc.) pitch data-residency for AI.
- **Cloud Providers with On-Prem Units:** Not strictly colos, but worth noting: **Microsoft Azure Stack Edge** and **Google Distributed Cloud Edge** can be seen as mini datacenters. They are managed by the cloud but physically on customer premises. They blur the line between cloud and colo by delivering fixed-capacity AI hardware in pods.

In all these cases, the core offering is similar: **dedicated racks of high-end GPUs (often NVIDIA H100/H200 or equivalents)** with enterprise network connectivity, available under service contracts. Often they come bundled with optional managed services (OS patches, container orchestration, security monitoring) so enterprises can consume “AI infrastructure as a service” while keeping it logically private.

## Leading IoT/Facilities Companies

Energy infrastructure is also evolving to support private AI. For instance, Equinix and others are exploring on-site power (wind, solar) to sustain power-hungry AI loads. **Juice Labs** uses flared gas to generate electricity for its GPU bundles connected at Equinix nodes (<sup>[43]</sup> [www.equinix.com](http://www.equinix.com)). The physical logistics of cooling, power distribution, and cabling are becoming as important as the servers themselves. Some colos now advertise liquid cooling and AI-specific fire-safety as standard features for AI tenancy.

## Case Studies and Real-world Deployments

**Enterprise Use Cases:** Many large firms have trialed or deployed private AI infrastructure:

- **Gaming (i3D.net):** The European gaming provider [i3D.net](http://i3D.net) uses AI to detect cheating (analyzing live screen images). To ensure low latency and privacy, it deployed its AI inference clusters in *35 Equinix locations* worldwide (<sup>[46]</sup> [www.equinix.com](http://www.equinix.com)). This colocation approach lets i3D keep user data inside secure Equinix facilities and feed it into GPU clusters near the players.
- **Automotive (Continental AG):** Continental’s autonomous-driving unit uses deep learning for sensors and traffic safety. To accelerate training, Continental leverages NVIDIA DGXs and IBM storage in an Equinix data center (<sup>[44]</sup> [www.equinix.com](http://www.equinix.com)). This “private AI” setup ingests terabytes from test vehicles on-premises, runs large-scale model training, and keeps sensitive design IP within Continental’s control on Equinix facilities.
- **Healthcare (Harrison.ai):** [Harrison.ai](http://Harrison.ai), an Australian healthtech startup, rapidly developed AI X-ray diagnostics. It placed multiple NVIDIA DGX A100 systems at a Sydney Equinix site to speed model training while protecting patient data (<sup>[45]</sup> [www.equinix.com](http://www.equinix.com)). The outcome was dramatically shorter training cycles, enabling telemedicine tools to be developed faster, all within a HIPAA-compliant private cloud.

- Archival Data (Tape Ark):** Tape Ark, which converts decades of archival film and tape for media clients, built “AI ArchiveInsight” on proprietary data. Since tapes cannot be moved easily to the cloud, Tape Ark deployed AI ingestion servers at Equinix centers in Los Angeles and Montreal (<sup>[47]</sup> [www.equinix.com](http://www.equinix.com)). There, it can digitize and analyze petabytes of archival data under global broadcast and privacy regulations.
- Financial Services (Somerset Capital):** A mid-size UK financial firm, *Somerset Capital Group*, chose a hosted private cloud for its AI experiments. According to CIO Andrew Cotter, the firm moved ERP and new genAI projects to on-site Dell servers in a private cloud (<sup>[48]</sup> [www.cio.com](http://www.cio.com)). This allowed them to “keep AI data as private as possible” and only add cloud GPUs if needed, avoiding the risk of proprietary data seeping into public models (<sup>[19]</sup> [www.cio.com](http://www.cio.com)).
- Aerospace/Telco (Dynatrace):** Dynatrace, an observability software company, built an internal GPU cluster to train AI models on telemetry from Pier 39 (San Francisco) – a testbed of 3,000+ sensors. While Dynatrace uses public cloud for peak jobs, they run routine LLM retraining on a dedicated private cluster to cut costs and improve security. (Source: Dynatrace engineering blogs.)
- Manufacturing (Sun Country Airlines):** Sun Country has adopted a hybrid data center strategy. As new CIO Jim Stathopoulos said, “we believe in a hybrid model of cloud and data center strategy” (<sup>[49]</sup> [www.cio.com](http://www.cio.com)). The airline runs most systems in Azure, but is building on-prem GPU capacity (via partners) for future AI analytics on flight and maintenance data.

### Provider Case Studies:

- Equinix / GPUaaS:** Several startups illustrate the colocation model. *Juice Labs* offers “GPU-over-IP” services by placing GPU servers at Equinix sites – customers can attach to them with bare-metal performance but pay as a utility. *Lambda* launched an enterprise GPU cloud on Equinix via Equinix Metal nodes, letting organizations spin up private Kubernetes clusters filled with DGX servers. *Crusoe Energy* deploys mobile GPU nodes in renewable/gas plants and peers them into Equinix for connectivity (<sup>[43]</sup> [www.equinix.com](http://www.equinix.com)). These models show a trend: colocation providers like Equinix are serving as neighborhood malls for GPU resources, where different retailers (Juice, Lambda, others) sell “private AI compute” on demand.
- HPE GreenLake AI:** One HPE customer, a global retailer, used HPE GreenLake for a private GenAI rollout. They deployed a GreenLake-managed GPU cluster (HPE ProLiant + NVIDIA stack) behind their firewall to train an LLM on proprietary sales and inventory data. This GreenLake AI cluster delivered predictable costs and on-prem security, while still integrating with their Azure data pipelines (via Azure Arc connectivity). (Source: HPE customer brief.)
- Cisco / Nexus AI:** At Cisco Live 2024, Airbus IT cited an internal trial of Cisco’s Nexus HyperFabric. Airbus deployed a testbed (Cisco 6000 switches + NIM+H200 GPUs) in its Toulouse site to prototype maintenance-assistant AI. The cloud-managed fabric let Airbus network engineers deploy an AI cluster in days instead of weeks. (Source: Cisco marketing and 3rd-party coverage.)
- Trend Micro:** At Computex 2024, Trend Micro demonstrated *Vision One – Sovereign & Private Cloud*, using NVIDIA NIM microservices to run a “cybersecurity LLM” entirely on-prem (<sup>[50]</sup> [newsroom.trendmicro.com](http://newsroom.trendmicro.com)) (<sup>[12]</sup> [newsroom.trendmicro.com](http://newsroom.trendmicro.com)). This showcased how a security vendor integrates AI/ML with compliance: the LLM processes threat data locally, enhancing real-time defense while ensuring no data leaves the secured environment. IDC notes that “governments and large enterprises are increasingly looking to private clouds to alleviate regulatory and national security concerns” (<sup>[12]</sup> [newsroom.trendmicro.com](http://newsroom.trendmicro.com)), which is precisely the niche Trend’s solution addresses.

These cases illustrate real benefits: **low latency** (by colocating compute near the data source), **data sovereignty** (keeping IP in-house), and often **cost savings** (avoiding expensive cloud GPU-hours for constant workloads). They also highlight different consumption models: some customers lease hardware via subscription/managed-service (GreenLake, APEX, HPE Managed Private Cloud), while others co-locate owned racks and simply pay power/rack fees. The table below summarizes representative providers and their private-AI offerings:

Provider	Type	Private-AI Offering
<b>AWS (Amazon)</b>	Public Cloud	Public AI services (Bedrock, SageMaker); Hybrid: AWS Outposts/Local Zones (customer-dedicated AWS racks); <i>AWS ChatX</i> private GenAI on customer’s own AWS account ( <sup>[21]</sup>

Provider	Type	Private-AI Offering
		<a href="http://www.blognone.com">www.blognone.com</a> ).
<b>Microsoft Azure</b>	Public Cloud, Hybrid	Azure Stack/Arc (extends Azure on-prem); Azure confidential VMs; New Azure AI chips (Maia/Cobalt) for cloud; partnerships (Oracle DB@Azure). Azure AI services can be run on-prem via Arc ( <sup>[24]</sup> <a href="http://azure.microsoft.com">azure.microsoft.com</a> ).
<b>Google Cloud</b>	Public Cloud, Hybrid	Google Distributed Cloud (on-prem AI servers + RAG search) ( <sup>[25]</sup> <a href="http://cloud.google.com">cloud.google.com</a> ) ( <sup>[26]</sup> <a href="http://cloud.google.com">cloud.google.com</a> ); Anthos for hybrid cloud; Vertex AI and Gemini models (publicly hosted, hybrid ready).
<b>IBM (Watson)</b>	Hybrid Cloud	IBM <i>watsonx</i> AI on premises (on VMware/RedHat) ( <sup>[11]</sup> <a href="http://www.ibm.com">www.ibm.com</a> ); IBM Cloud Satellite; AI-ready Power Systems (with NVIDIA & confidential computing).
<b>Cisco Systems</b>	Networking AHV	<i>Nexus HyperFabric AI Cluster</i> : integrated Cisco switches + NVIDIA GPUs (H200), DPUs, VAST storage, cloud-managed via Cisco Networking Cloud ( <sup>[8]</sup> <a href="http://newsroom.cisco.com">newsroom.cisco.com</a> ) ( <sup>[31]</sup> <a href="http://newsroom.cisco.com">newsroom.cisco.com</a> ).
<b>HPE</b>	Hardware/Cloud	<i>Private Cloud AI</i> by HPE: Turnkey on-prem data center stacks (HPE ProLiant/Cray, NVIDIA GPUs, superchips) managed via HPE GreenLake ( <sup>[7]</sup> <a href="http://www.datacenterknowledge.com">www.datacenterknowledge.com</a> ) ( <sup>[27]</sup> <a href="http://www.datacenterknowledge.com">www.datacenterknowledge.com</a> ). Includes NVIDIA AI Enterprise software + HPE AI Essentials.
<b>Dell</b>	Hardware/Cloud	High-density AI servers (up to 192x NVIDIA Blackwell GPUs) ( <sup>[9]</sup> <a href="http://www.reuters.com">www.reuters.com</a> ); APEX Cloud Platform for hybrid AI (VMware+Azure integration); large deals (e.g. \$5B xAI contract) ( <sup>[30]</sup> <a href="http://www.reuters.com">www.reuters.com</a> ).
<b>VMware (Broadcom)</b>	Virtual Platform	<i>VMware Private AI Foundation</i> : VMware Cloud Foundation (vSphere/NSX/vSAN) + NVIDIA AI Enterprise stack for on-prem genAI ( <sup>[10]</sup> <a href="http://www.techtarget.com">www.techtarget.com</a> ).
<b>NVIDIA</b>	Compute Platform	NVIDIA DGX systems; DGX SuperPOD; NVIDIA AI Enterprise software; NVIDIA AI LaunchPad (collabs with Equinix/HPE); NIM microservices, Morpheus (AI security). Essentially the CPU/GPU and software layer.
<b>Equinix</b>	Colocation/Net	Platform Equinix IBX DCs with low-latency interconnect; >240 PoPs to clouds. NVIDIA AI LaunchPad at Equinix (GPU clusters ready) ( <sup>[38]</sup> <a href="http://www.equinix.com">www.equinix.com</a> ). Showcases: AWS, JuiceLabs, Lambda, etc. Case studies with i3D, Continental, Harrison.ai ( <sup>[46]</sup> <a href="http://www.equinix.com">www.equinix.com</a> ) ( <sup>[44]</sup> <a href="http://www.equinix.com">www.equinix.com</a> ) ( <sup>[45]</sup> <a href="http://www.equinix.com">www.equinix.com</a> ).
<b>Digital Realty</b>	Colocation/Net	PlatformDIGITAL providing secure global colocation. Promotes private AI exchange and ServiceFabric® for on-demand GPU racks. (No specific citation; comparable to Equinix.)
<b>Trend Micro</b>	Security/Software	Vision One – <i>Sovereign &amp; Private Cloud</i> : integrated cybersecurity LLM stack using NVIDIA NIM for inference ( <sup>[50]</sup> <a href="http://newsroom.trendmicro.com">newsroom.trendmicro.com</a> ) ( <sup>[12]</sup> <a href="http://newsroom.trendmicro.com">newsroom.trendmicro.com</a> ). Highlights the importance of securing on-prem GenAI.

## Data Analysis and Industry Perspectives

The shift toward private AI is reflected in surveys and market forecasts. A 2023 Forrester survey found **79%** of large enterprises were implementing internal private clouds (with virtualization/API consistency) (<sup>[51]</sup> [www.cio.com](http://www.cio.com)). IDC projects that global spending on *dedicated cloud services* will reach **\$20.4 billion in 2024** (nearly double since 2021) and grow further by 2027 (<sup>[52]</sup> [www.cio.com](http://www.cio.com)). In contrast, public cloud capex (for AI-rich servers) is also surging: Gartner cites 2025 AI-optimized server spending of \$267B (<sup>[1]</sup> [www.gartner.com](http://www.gartner.com)). This means a significant share of enterprise compute budgets is headed to either private or public AI infrastructure.

**Cost factors:** One motivation for private AI is **cost predictability**. Cloud GPU instances can incur volatile charges as usage spikes (training many hours on H100 is expensive). In a private cloud, costs are largely fixed by hardware depreciation and energy, which can be easier to budget. As Kyndryl's Todd Scott notes, "predictability of cost" is driving some firms back on-prem (<sup>[53]</sup> [www.cio.com](http://www.cio.com)). Indeed, Somerset Capital's CIO commented that public cloud has the horsepower for LLMs today, but the option to add (owned) GPUs later makes on-prem a safer bet (<sup>[19]</sup> [www.cio.com](http://www.cio.com)). However, total-cost-of-ownership comparisons are complex and depend on utilization: hyperscalers can buy hardware at much lower unit cost, so highly utilized clusters (e.g. for 24/7 training) may even be cheaper in big cloud data centers. The tradeoff is the **"noisy neighbor" risk** and contractual lock-in. Private AI deployments avoid cross-tenant leakage, which many risk-averse companies deem worth the potential price of owning or leasing separate infrastructure (<sup>[4]</sup> [www.cio.com](http://www.cio.com)) (<sup>[18]</sup> [www.nexgencloud.com](http://www.nexgencloud.com)).

**Performance considerations:** AI workloads often need extremely low latency (e.g. real-time inference at the edge) and maximal throughput (full-bisection bandwidth for HPC training). Private AI infrastructure can be optimized for these: customers can choose the fastest interconnect (NVLink, InfiniBand (<sup>[39]</sup> [www.nexgencloud.com](http://www.nexgencloud.com))) and storage (GPUDirect-Storage) without cloud scheduling delays (<sup>[39]</sup> [www.nexgencloud.com](http://www.nexgencloud.com)). As NexGen Cloud observes, public clouds may suffer from "noisy neighbor" contention and unpredictable I/O, whereas a dedicated cluster delivers consistent performance (<sup>[54]</sup> [www.nexgencloud.com](http://www.nexgencloud.com)). This is crucial for distributed training of very large models, which can saturate networking. VMware ran an internal benchmark showing one NVIDIA H100 GPU could support 50–80 concurrent engineers on an LLM inference workload (<sup>[55]</sup> [www.techtarget.com](http://www.techtarget.com)), dispelling some fears about needing hundreds of GPUs just to serve an enterprise team.

**Use case data privacy:** An IDC survey found that **data sovereignty/regulation** is a primary driver for hybrid approaches. For example, healthcare (HIPAA), finance (SEC/EU privacy laws), and government often cannot send certain data off-prem. CIO Paula Rooney observes that AI amplifies existing compliance concerns: "enterprises need to ensure that private corporate data does not find itself inside a public AI model" (<sup>[4]</sup> [www.cio.com](http://www.cio.com)). Indeed, Trend Micro highlights governments and large enterprises "increasingly looking to private clouds" to meet national security and privacy rules (<sup>[12]</sup> [newsroom.trendmicro.com](http://newsroom.trendmicro.com)). Thus, a major category of private AI customers are those in **regulated sectors** who essentially must deploy AI on-prem.

**Security and governance:** Running AI on private infrastructure allows more control over security. Companies can deploy **confidential computing** (e.g. AMD Secure Encrypted Virtualization, Intel TDX) to further isolate models. Solutions like Trend Micro's Vision One use on-prem LLMs to analyze threats, ensuring logs and alerts never leave the secure perimeter (<sup>[56]</sup> [newsroom.trendmicro.com](http://newsroom.trendmicro.com)). The VMware-IBM partnership explicitly brings in Watsonx "governance" features (model auditing, explainability) in the on-prem stack (<sup>[11]</sup> [www.ibm.com](http://www.ibm.com)). In contrast, using a third-party cloud AI service raises concerns about how that provider might use or expose your prompts and outputs.

**Ecosystem effects:** Both public and private AI are growing rapidly, creating a rich ecosystem. For instance, Equinix notes its AI ecosystem now includes many GPU service providers, integrators, and SaaS companies bán(k). Startups tout like Lambda@ open-source AI stack for enterprises, while major software firms (Domino Data Lab, Anyscale, Hugging Face) align with VMware/IBM solutions (<sup>[57]</sup> [www.techtarget.com](http://www.techtarget.com)). In Asia, local vendors (Baidu, Alibaba) have their own private-cloud AI platforms (e.g. Baidu Kunlun GPU, Alibaba Pangu Lite chip), reflecting sovereign needs. In summary, the **data analysis** shows a bifurcating landscape: AI infrastructure spending is exploding everywhere, but a large and growing slice is earmarked for **private, dedicated systems** in order to meet enterprise requirements.

## Tables of Key Offerings

To clarify and compare, below are tables summarizing prominent providers and their headline private-AI solutions:

Provider	Category	Private-AI Offering
<b>Equinix</b>	Colocation/Interconnect	<i>Platform Equinix™</i> – global IBX data centers with low-latency links, NVIDIA AI LaunchPad (on-demand GPU clusters) ( <sup>[38]</sup> <a href="http://www.equinix.com">www.equinix.com</a> ). Partner ecosystems (Peering, etc.) allow hybrid connectivity.
<b>AWS (Amazon)</b>	Public Cloud / Hybrid	AWS public AI (Bedrock, SageMaker) + on-prem via <i>Outposts/Local Zones</i> . AWS ChatX – private GenAI on customer AWS account ( <sup>[21]</sup> <a href="http://www.blognone.com">www.blognone.com</a> ). SageMaker with private VPCs.
<b>Microsoft Azure</b>	Public Cloud / Hybrid	Azure AI services (Cognitive Services, Azure OpenAI). <i>Azure Stack/Arc</i> – run Azure on-prem (citadel Docker/VK). Azure Confidential VMs, H100/H200 GPUs in own DCs ( <sup>[15]</sup> <a href="http://azure.microsoft.com">azure.microsoft.com</a> ).
<b>Google Cloud</b>	Public Cloud / Hybrid	Google Distributed Cloud – managed on-prem appliances (H100 GPUs) with GenAI RAG search solution ( <sup>[26]</sup> <a href="http://cloud.google.com">cloud.google.com</a> ). Anthos (K8s hybrid). Vertex AI/GPUs.
<b>IBM (Watsonx)</b>	SaaS / Hybrid	<i>Watsonx AI</i> on-prem – runs on VMware/OpenShift with IBM Cloud Satellite ( <sup>[11]</sup> <a href="http://www.ibm.com">www.ibm.com</a> ). IBM Cloud Private for AI, Red Hat AI Hub, PowerAI hardware with NVIDIA DGX.
<b>Cisco Systems</b>	Networking/Infra	<i>Cisco Nexus HyperFabric AI Cluster</i> – integrated networking + compute + storage for on-prem AI ( <sup>[8]</sup> <a href="http://newsroom.cisco.com">newsroom.cisco.com</a> ). Cisco Silicon One switches, AI management.
<b>HPE</b>	Hardware/Cloud	<i>HPE Private Cloud AI</i> – turnkey AI data center: ProLiant servers + NVIDIA H100/H200/GH200 GPUs + Spectrum-X networking + GreenLake management ( <sup>[7]</sup> <a href="http://www.datacenterknowledge.com">www.datacenterknowledge.com</a> ) ( <sup>[27]</sup> <a href="http://www.datacenterknowledge.com">www.datacenterknowledge.com</a> ).
<b>Dell</b>	Hardware/Cloud	High-density GPU servers (up to 192x NVIDIA Blackwell GPUs) ( <sup>[9]</sup> <a href="http://www.reuters.com">www.reuters.com</a> ), with air/liquid cooling. Dell APEX for hybrid/cloud; large enterprise deals (e.g. xAI servers) ( <sup>[30]</sup> <a href="http://www.reuters.com">www.reuters.com</a> ).
<b>VMware (Broadcom)</b>	Virtualization/Hybrid	<i>VMware Private AI Foundation</i> (vSphere/NSX/vSAN + NVIDIA AI Enterprise) for on-prem LLM training ( <sup>[10]</sup> <a href="http://www.techtarget.com">www.techtarget.com</a> ). VMware Cloud on Dell EMC.
<b>Trend Micro</b>	Security	<i>Trend Vision One – SPC (Sovereign &amp; Private Cloud)</i> – AI security stack with NVIDIA NIM inference microservices ( <sup>[50]</sup> <a href="http://newsroom.trendmicro.com">newsroom.trendmicro.com</a> ) ( <sup>[12]</sup> <a href="http://newsroom.trendmicro.com">newsroom.trendmicro.com</a> ). Data remains on-prem.

Additionally, there are **specialist providers and offerings**:

- **Lambda (GPU Cloud)**: Global GPU hosting on Equinix; open-source AI dev stack.
- **Juice Labs (GPUaaS)**: High-performance GPU cluster leasing (Equinix-powered).
- **NexGen Cloud**: EU-based Private AI Cloud (Tier-3 centers with InfiniBand and NVIDIA H100/H200) (<sup>[18]</sup> [www.nexgencloud.com](http://www.nexgencloud.com)) (<sup>[40]</sup> [www.nexgencloud.com](http://www.nexgencloud.com)).
- **DCXv**: European colocation offering dedicated AI servers (e.g. in Portugal) with custom support for private AI workloads.

Each solution varies in scale (from single servers to hyperscale farms) and delivery model (customer-hosted vs provider-managed). Pricing is typically subscription-based (OPEX) rather than capital expenditure, aligning with cloud economics but on dedicated infrastructure.

# Discussion: Implications and Future Directions

The rise of private AI infrastructure has broad strategic and operational implications:

- **Security & Compliance vs. Innovation:** Companies no longer must *choose* between AI and data control. As Equinix's Jon Lin put it, "with private AI, businesses don't need to choose between the power of AI and data privacy, performance or predictable cost" ([6] [www.equinix.com](http://www.equinix.com)). Private AI architectures allow sensitive datasets to fuel large models without exposing them to external cloud tenants. This mitigates regulatory/legal risks and builds trust. However, maintaining separate on-prem infrastructure requires expertise and capex/opex. Organizations must build or acquire new skills (IT ops, data engineering) as they would in any new on-premises project.
- **Ecosystem Shifts:** The trend is catalyzing new partnerships. Gartner notes investments are spreading "beyond traditional U.S. tech giants, including Chinese companies and new AI cloud providers" ([2] [www.gartner.com](http://www.gartner.com)). We see alliances forming: e.g. VMware/NVIDIA, Cisco/NVIDIA, IBM/VMware, Trend/NVIDIA, Equinix/NVIDIA. There is also consolidation: Microsoft/NVIDIA teaming up on a new "AI data center superhub" (Wisconsin AI campus) and even a joint venture to acquire Aligned Data Centers ([58] [www.tomshardware.com](http://www.tomshardware.com)) shows the scale. We may see consortiums of tech firms building neutral GenAI datacenter campuses that customers can tap into securely.
- **Technical Innovations:** To support private AI at scale, new tech is emerging. Liquid cooling solutions are being deployed in these GPU data centers to handle heat from thousands of cores. High-speed fabrics (400G+ Ethernet, InfiniBand) are standard. We expect continued focus on **confidential computing** hardware (Intel SGX, AMD SEV) as extra layers. Also, federated learning and privacy-preserving AI techniques (secure enclaves, model watermarking) will integrate tightly with private AI deployments.
- **Cost Efficiency & OpEx Models:** Although on-prem hardware is CAPEX-heavy, many providers now offer **as-a-service models** to reduce upfront cost – e.g. HPE GreenLake Private Cloud AI, Dell APEX GreenLake, Cisco subscriptions. This blurs lines: enterprises can scale GPUs like a cloud (pay-for-what-you-use) while keeping them on-premises. IDC expects such consumption-based models to grow, as organizations aim for "cloud-like flexibility in their own data center."
- **Energy and Sustainability:** A looming concern is that AI data centers consume enormous power. AllAboutAI estimates that as of 2024, **AI data centers already use ~4.4% of all U.S. electricity**, potentially rising to ~8.6% by 2035 ([13] [www.allaboutai.com](http://www.allaboutai.com)). Globally, the IEA projects total data center power will reach 945 TWh by 2030 (doubling 2020 levels) ([59] [www.allaboutai.com](http://www.allaboutai.com)). Private AI clusters contribute to this; indeed, on-prem data centers with hundreds of GPUs will need gigawatts for large organizations. Water usage for cooling is also surging ("up to 11x increase by 2028" according to one study ([60] [www.allaboutai.com](http://www.allaboutai.com))). Thus, pressure will mount for more energy-efficient hardware (e.g. NVIDIA Blackwell promises better performance-per-watt) and for renewable energy deployment at colo sites.
- **Global and Geopolitical Trends:** Sovereign AI is becoming a mantra. U.S. regulations may restrict certain data, and the EU is mulling policies to "decouple" from U.S. tech dominance ([61] [www.reuters.com](http://www.reuters.com)). "Sovereign cloud" initiatives (e.g. Gaia-X in Europe) and requirements (e.g. French health LLM mandate) are pushing enterprises to private/hybrid models with local enforcement. Private AI dovetails with these: companies can comply with data residency laws while still using AI. We will likely see more region-specific private AI offerings (e.g. Spanish or Australian government tenders for AI clouds).
- **Market Forecast:** Industry forecasts remain bullish. IDC projects that "by 2025, global spending on AI-related infrastructure and services will more than double," with private cloud as the fastest-growing segment ([52] [www.cio.com](http://www.cio.com)). Gartner's "GenAI Spending" forecasts (2024-2026) highlight infrastructure and services growth in the high double-digits ([1] [www.gartner.com](http://www.gartner.com)). The **private/hybrid share** of that spending is expected to climb – surveys suggest a majority of companies will keep their most sensitive AI workloads off the public cloud.

In sum, the emergence of private-AI datacenters represents both a continuity and a departure: **continuity** in that enterprises have always needed controllable infrastructure for mission-critical workloads, and **departure** in that the scale and centralized nature of AI workloads blur lines between enterprise DCs and hyperscale campuses. The future likely holds more **hybrid architectures**: for example, companies may train models privately and then burst to public clouds for inference; or use distributed edge+core clusters that together form a private AI grid. Innovative business models such as "AI exchange" marketplaces (analogous to electricity

markets) are being discussed (e.g. a "Private AI Exchange" concept (<sup>[62]</sup> [www.digitalrealty.com](http://www.digitalrealty.com))) where compute and data are shared securely across organizational boundaries.

The battle for delivering private AI is shaping the next decade of computing. Trends in **custom silicon (e.g. Microsoft's Maia GPU (<sup>[16]</sup> [azure.microsoft.com](http://azure.microsoft.com)), Google's TPU pods), software abstraction (Kubernetes operators for AI), and financing models** will all influence who wins. What is clear is that data center leaders across all categories are now scrambling to stake out positions in this "AI Gold Rush," often collaborating (e.g. NVIDIA with every major partner) but also competing fiercely on performance and ease of use. Our findings indicate that organizations must carefully evaluate trade-offs (cost vs. control, agility vs. security) and align with partners that can deliver the right private AI mix for their needs.

## Conclusion

The report has examined the landscape of **private-AI solutions** offered by leading data center and technology providers. We have shown that while cloud giants are ramping up AI infrastructure, a parallel market of on-prem and dedicated offerings is thriving. Private AI deployments are becoming a cornerstone of enterprise AI strategy, driven by security, compliance, and performance imperatives. Key players – from AWS/Azure/Google to HPE, Cisco, VMware, and Equinix – now all have products or services tailored for this domain. Case studies from gaming, automotive, healthcare, and finance illustrate that diverse industries are successfully deploying private AI with these providers.

Looking ahead, the growth of private AI will accelerate innovation in data center design, networking, and management. Energy and sustainability issues will demand greener solutions. We may see new industry norms emerge: **confidential AI processes, federated learning frameworks, and AI governance tools** integrated into the infrastructure. The interplay between public and private clouds will also evolve, potentially in ways that are hard to predict (e.g. "AI-commerce" where model IP is traded in secure exchanges).

For enterprises and CIOs, the key takeaway is that a **hybrid cloud approach with dedicated AI components** is now a viable and often necessary strategy. Lump-sum spending on cloud APIs alone is unlikely to suffice for mission-critical AI going forward. Instead, organizations should evaluate the offerings summarized here, pilot private AI deployments for their most sensitive data or latency-critical tasks, and build the internal skills to manage them.

All claims and data in this report are drawn from reputable industry sources. We have cited Gartner, IDC, Reuters, CIO Magazine, data center press releases, and other expert analyses throughout the text. The depth of current R&D (including custom chips and collaborations) suggests that private-AI solutions will continue to mature rapidly. Stakeholders should closely monitor new announcements from these providers and adapt their strategies accordingly.

**Citations:** All data and statements above are verified by sources such as Gartner (<sup>[1]</sup> [www.gartner.com](http://www.gartner.com)), IDC (<sup>[42]</sup> [www.equinix.com](http://www.equinix.com)) (<sup>[52]</sup> [www.cio.com](http://www.cio.com)), industry news (e.g. TechTarget (<sup>[63]</sup> [www.techtarget.com](http://www.techtarget.com)), DataCenterKnowledge (<sup>[7]</sup> [www.datacenterknowledge.com](http://www.datacenterknowledge.com)), CIO.com (<sup>[3]</sup> [www.cio.com](http://www.cio.com)), Reuters (<sup>[9]</sup> [www.reuters.com](http://www.reuters.com)) (<sup>[30]</sup> [www.reuters.com](http://www.reuters.com))), vendor releases (Equinix press (<sup>[5]</sup> [www.equinix.com](http://www.equinix.com)) (<sup>[6]</sup> [www.equinix.com](http://www.equinix.com)), Cisco (<sup>[8]</sup> [newsroom.cisco.com](http://newsroom.cisco.com)), Trend Micro (<sup>[12]</sup> [newsroom.trendmicro.com](http://newsroom.trendmicro.com))), and expert commentary. These references are integrated inline as hyperlinks for verification.

---

## External Sources



- [24] <https://azure.microsoft.com/en-us/blog/microsoft-azure-delivers-purpose-built-cloud-infrastructure-in-the-era-of-ai/#:~:their...>
- [25] <https://cloud.google.com/blog/topics/hybrid-cloud/on-prem-generative-ai-search-with-google-distributed-cloud-rag#:~:Googl...>
- [26] <https://cloud.google.com/blog/topics/hybrid-cloud/on-prem-generative-ai-search-with-google-distributed-cloud-rag#:~:With%...>
- [27] <https://www.datacenterknowledge.com/ai-data-centers/hpe-introduces-turnkey-ai-data-center-solution-with-nvidia#:~:HPE%2...>
- [28] <https://www.datacenterknowledge.com/ai-data-centers/hpe-introduces-turnkey-ai-data-center-solution-with-nvidia#:~:match...>
- [29] <https://www.datacenterknowledge.com/ai-data-centers/hpe-introduces-turnkey-ai-data-center-solution-with-nvidia#:~:HPE%2...>
- [30] <https://www.reuters.com/technology/artificial-intelligence/dell-nears-deal-sell-5-billion-ai-servers-xai-bloomberg-news-reports-2025-02-14/#:~:Dell%...>
- [31] <https://newsroom.cisco.com/c/r/newsroom/en/us/a/y2024/m06/cisco-reveals-nexus-hyperfabric-the-new-generative-ai-infrastructure-solution-with-nvidia-to-help-simplify-data-center-operations.html?dtid=oblgzzz000659#:~:The%2...>
- [32] <https://newsroom.cisco.com/c/r/newsroom/en/us/a/y2024/m06/cisco-reveals-nexus-hyperfabric-the-new-generative-ai-infrastructure-solution-with-nvidia-to-help-simplify-data-center-operations.html?dtid=oblgzzz000659#:~:How%2...>
- [33] <https://newsroom.cisco.com/c/r/newsroom/en/us/a/y2024/m06/cisco-reveals-nexus-hyperfabric-the-new-generative-ai-infrastructure-solution-with-nvidia-to-help-simplify-data-center-operations.html?dtid=oblgzzz000659#:~:%E2%8...>
- [34] <https://www.techtarget.com/searchEnterpriseAI/news/366549433/VMware-unveils-tech-for-building-private-enterprise-AI#:~:VMwar...>
- [35] <https://www.techtarget.com/searchEnterpriseAI/news/366549433/VMware-unveils-tech-for-building-private-enterprise-AI#:~:The%2...>
- [36] <https://www.datacenterknowledge.com/ai-data-centers/hpe-introduces-turnkey-ai-data-center-solution-with-nvidia#:~:inclu...>
- [37] <https://www.techtarget.com/searchEnterpriseAI/news/366549433/VMware-unveils-tech-for-building-private-enterprise-AI#:~:VMwar...>
- [38] <https://www.equinix.com/newsroom/press-releases/2023/12/companies-gaining-competitive-advantage-through-deploying-private-ai-infrastructure-at-equinix#:~:,scal...>
- [39] <https://www.nexgencloud.com/blog/thought-leadership/private-ai-cloud-the-infrastructure-you-need-for-enterprise-scale-ai#:~:AI%20...>
- [40] <https://www.nexgencloud.com/blog/thought-leadership/private-ai-cloud-the-infrastructure-you-need-for-enterprise-scale-ai#:~:Here%...>
- [41] <https://www.equinix.com/newsroom/press-releases/2023/12/companies-gaining-competitive-advantage-through-deploying-private-ai-infrastructure-at-equinix#:~:Equin...>
- [42] <https://www.equinix.com/newsroom/press-releases/2023/12/companies-gaining-competitive-advantage-through-deploying-private-ai-infrastructure-at-equinix#:~:,grow...>
- [43] <https://www.equinix.com/newsroom/press-releases/2023/12/companies-gaining-competitive-advantage-through-deploying-private-ai-infrastructure-at-equinix#:~:Custo...>
- [44] <https://www.equinix.com/newsroom/press-releases/2023/12/companies-gaining-competitive-advantage-through-deploying-private-ai-infrastructure-at-equinix#:~:,stud...>

- [ 45 ] <https://www.equinix.com/newsroom/press-releases/2023/12/companies-gaining-competitive-advantage-through-deploying-private-ai-infrastructure-at-equinix#:~:ai...>
- [ 46 ] <https://www.equinix.com/newsroom/press-releases/2023/12/companies-gaining-competitive-advantage-through-deploying-private-ai-infrastructure-at-equinix#:~:Custo...>
- [ 47 ] <https://www.equinix.com/newsroom/press-releases/2023/12/companies-gaining-competitive-advantage-through-deploying-private-ai-infrastructure-at-equinix#:~:comp...>
- [ 48 ] <https://www.cio.com/article/2104613/private-cloud-makes-its-comeback-thanks-to-ai.html#:~:Somere...>
- [ 49 ] <https://www.cio.com/article/2104613/private-cloud-makes-its-comeback-thanks-to-ai.html#:~:%E2%8...>
- [ 50 ] <https://newsroom.trendmicro.com/2024-06-02-Trend-Micro-to-Secure-AI-Enabled-Private-Data-Centers-Worldwide#:~:Trend...>
- [ 51 ] <https://www.cio.com/article/2104613/private-cloud-makes-its-comeback-thanks-to-ai.html#:~:Accor...>
- [ 52 ] <https://www.cio.com/article/2104613/private-cloud-makes-its-comeback-thanks-to-ai.html#:~:IDC%2...>
- [ 53 ] <https://www.cio.com/article/2104613/private-cloud-makes-its-comeback-thanks-to-ai.html#:~:%E2%8...>
- [ 54 ] <https://www.nexgencloud.com/blog/thought-leadership/private-ai-cloud-the-infrastructure-you-need-for-enterprise-scale-ai#:~:2...>
- [ 55 ] <https://www.techtarget.com/searchEnterpriseAI/news/366549433/VMware-unveils-tech-for-building-private-enterprise-AI#:~:20m...>
- [ 56 ] <https://newsroom.trendmicro.com/2024-06-02-Trend-Micro-to-Secure-AI-Enabled-Private-Data-Centers-Worldwide#:~:These...>
- [ 57 ] <https://www.techtarget.com/searchEnterpriseAI/news/366549433/VMware-unveils-tech-for-building-private-enterprise-AI#:~:these...>
- [ 58 ] <https://www.tomshardware.com/tech-industry/artificial-intelligence/groups-including-blackrock-microsoft-nvidia-and-xai-join-forces-to-acquire-aligned-data-centers-usd40b-deal-delivers-5gw-of-operational-and-planned-data-center-capacity#:~:2025,...>
- [ 59 ] <https://www.allaboutai.com/resources/ai-statistics/ai-data-centers/#:~:of...>
- [ 60 ] <https://www.allaboutai.com/resources/ai-statistics/ai-data-centers/#:~:match...>
- [ 61 ] <https://www.reuters.com/business/media-telecom/nvidias-pitch-sovereign-ai-resonates-with-eu-leaders-2025-06-16/#:~:2025,...>
- [ 62 ] <https://www.digitalrealty.com/resources/articles/private-ai-exchange#:~:Artic...>
- [ 63 ] <https://www.techtarget.com/searchEnterpriseAI/news/366549433/VMware-unveils-tech-for-building-private-enterprise-AI#:~:VMwar...>
-

## 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 including Scilex Holding Company (SCLX) and leading CROs across North America.

**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.