

WHITE PAPER

## Rearchitecting data centers for AI workloads

According to some estimates, over 80% of AI projects<sup>1</sup> fail to achieve all their goals. There are several common reasons for this, one of them being a lack of adequate data center infrastructure.

While most IT leaders realize AI can't happen without cutting-edge infrastructure, building and optimizing an AI-ready data center is undeniably complex.

At the heart of this challenge is that AI isn't a single process but a set of processes happening in different places, each with its own technical requirements. Before you can train an AI model, you need to collect, process and store vast amounts of data, often across a multitude of systems.

Training then requires specialized hardware, such as high-end GPUs. Once trained, you need the low-latency, high-efficiency CPUs to generate real-time

insights, all while moving data seamlessly between storage, compute and edge environments.

Each workload places unique demands on the data center. As such, optimizing data center infrastructure to support AI projects requires a fundamentally different approach compared with more traditional enterprise workloads, such as virtualization or delivery of cloud applications.

This paper explores strategies for building an efficient data center that can accommodate the ever-growing demands of AI workloads.

Produced by

  
TechTarget

## Strategy 1: Prepare for AI with high-throughput data collection and storage

AI is tremendously data-hungry. The data sets used to train large mainstream models such as ChatGPT-4 or Google Gemini run into dozens of terabytes. Even smaller models, designed for specific use cases, require about 10 times the amount of data than the number of parameters, or the number of variables the model learns during training. ChatGPT-4, for example, has an estimated 1.8 trillion parameters.<sup>2</sup> Even the smallest AI models have millions of parameters.

Training data can significantly burden storage infrastructure, particularly as AI models grow in scale and complexity. Moreover, training data must often be ingested from a wide range of sources, such as internal databases, IoT devices, public data sets and synthetic data — all of which needs to be stored, organized and made readily available for model training.

Once collected, training data must be stored so it can be efficiently retrieved and processed, while also supporting quality control and governance standards.

### Key challenges in AI data collection and storage

- Data quality: Legacy data center architecture often lacks the tools to detect and resolve issues such as duplication, missing values or labeling errors.
- Data structure: Older systems may struggle to efficiently organize, tag or standardize unstructured data, thus reducing training efficiency.
- Volume and complexity: Traditional data centers are rarely equipped to handle the sheer size and complexity of the data required for training AI models.
- Performance bottlenecks: Existing storage solutions, particularly conventional hard drives, lack the necessary throughput, resulting in major data-processing delays.
- Network performance: General-purpose data centers often lack the high-bandwidth interfaces needed to efficiently handle the continuous transfer of huge volumes of data.



## **Step 1: Choose scalable, high-performance storage mediums**

To meet data collection and storage demands, the AI-ready data center must use modernized storage architectures that minimize bottlenecks and support high-speed data retrieval.

While hard drives remain a practical choice for low-cost archival storage, they can't meet the performance demands of real-time data access. An optimal production environment should instead use nonvolatile memory express solid state drives (NVMe SSDs). These offer lower latency and higher throughput, with data-transfer speeds up to 35 times higher than conventional hard drives and about 7 times higher than standard SSDs.

Scalability is another key factor since AI projects require continuous ingestion of new data for training future iterations. Solutions such as data lakes and object-based storage can support increasingly massive and diverse data sets without degrading performance, ensuring data ingestion and retrieval remain efficient — even as AI workloads become more complex.

## **Step 2. Optimize compute infrastructure for AI data pipelines**

Beyond selecting the right storage medium, AI-ready data centers require compute infrastructure optimized for data ingestion and preprocessing, not just training itself. After all, training data must be efficiently structured, validated and organized before training can begin.

Server-grade processors are essential for streamlining these processes. High-core-count architectures, with high memory bandwidth and caching, support rapid indexing and metadata tagging across AI data sets. Additionally, support for parallelized ingestion workloads enables data centers to process multiple streams of incoming data simultaneously, thus reducing bottlenecks in AI workflows.

Memory bandwidth is another key factor in ensuring efficient data movement. Solutions such as DDR5 and PCIe 5.0 can enhance transfer speeds between storage and compute environments, thereby reducing preprocessing times and ensuring that data is readily accessible for training.

## **Step 3. Build a high-performance network for AI data movement**

The AI-ready data center must also have the right networking architecture to support fast, congestion-free data movement between acquisition points, storage systems and processing units. Without such an environment, bottlenecks will emerge that can slow AI workloads before training even begins.

Data centers require networking that can handle data transmission at massive scale with minimal latency. For instance, smart data-processing units (DPUs) — themselves augmented by machine learning — automatically offload network-related tasks from the CPU to reduce processing overhead while accelerating data flow.

Moreover, high-bandwidth interconnects, such as PCIe and Compute Express Link (CXL) enable rapid data transfer between storage and compute environments for optimal AI data retrieval.



## Strategy 2: Maximizing compute power for high-performance model training

Model training is easily the most resource-intensive phase of the AI lifecycle. Computational demands are enormous, making raw compute power a common bottleneck in training large-scale AI models.

Not only does data center infrastructure need to accommodate rapid access to vast data sets, but it also needs high-performance computing (HPC) capabilities to handle the myriad complex computations inherent in training any AI model.

That said, the computational requirements for AI model training vary greatly, depending on the size and complexity of the model. For example, a small-scale model may have only a few billion parameters, in which case, a single high-end GPU might be enough to handle training over several days. However, for larger-scale models, such as general-purpose LLMs, training may span server clusters equipped with dozens of GPUs working around the clock for several weeks. To give a real-world example, ChatGPT 4.0 was trained over several days using about 25,000 GPUs.<sup>3</sup>

### Key challenges in AI model training and computation

- Computational demands: General-purpose data centers primarily use CPUs to handle a broad range of tasks, but they lack the parallel processing required for efficient training at scale.
- High energy consumption: Many existing data centers lack the required power. Even training smaller models with just a few billion parameters can consume thousands of kilowatt-hours — orders of magnitude more than traditional enterprise workloads.
- Heat generation and cooling: Conventional cooling typically can't keep up with the heat generated by high-performance components working 24/7.

### Step 1: Reduce computational bottlenecks with AI-ready hardware

General-purpose data centers are primarily designed for the sequential processing often used in many traditional workloads, but they fall short in large-scale parallel-processing workloads, such as AI model training.

Since training involves processing millions of computations simultaneously, it's impractical to rely on mainstream CPUs alone. For maximum efficiency and shorter training cycles, AI workloads can benefit from an architecture where high-performance CPUs handle orchestration, data preprocessing and memory management, while GPUs or AI accelerators handle parallelized model training workflows.

AI-specialized hardware, such as next-generation GPUs or AI accelerators, contains thousands of cores, as opposed to the triple-digit core counts in the leading server-grade CPUs. However, these various processing units all complement one another in the AI-ready data center since they excel at handling different tasks.





## **Step 2: Address energy-consumption demands in the AI data center**

AI model training consumes massive amounts of energy, with larger models requiring several orders of magnitude more energy than small, single-purpose models with just a few billion parameters. Also, AI training involves sustained, resource-hungry workloads, unlike the short bursts of computation used for inference. Standard CPUs aren't energy-efficient enough, leading to wasted energy consumption over excessively long training times.

AI-optimized processors have lower power-to-computation ratios and can maintain the high compute demands of AI model training without sacrificing performance or drawing excessive amounts of energy. With the right software and orchestration layer, the AI data center can dynamically balance workloads to ensure they run on the most efficient compute nodes, while more energy-efficient interconnects can reduce wasted compute cycles.

## **Step 3: Manage heat output with advanced cooling**

Naturally, the long compute cycles of AI model training generate a lot of heat, even with the latest and most efficient AI accelerators, GPUs and CPUs. Power density in the AI data center tends to be enormous, with AI-ready hardware packing extreme power in ever-smaller footprints. Inadequate cooling can drastically reduce the life span of components and result in soaring energy use, necessitating advanced thermal management and next-generation cooling.

Direct-to-chip liquid cooling is far more energy-efficient and effective than air cooling, especially in high-density GPU clusters. For even more extreme heat loads, immersion cooling, where servers are submersed in nonconductive cooling liquids, is emerging as a promising alternative, especially in exascale supercomputers. Finally, the software layer can use machine learning algorithms to monitor and dynamically adjust cooling settings, such as fan speeds, in real time.

## Strategy 3: Optimizing networking and compute for real-time inference

After training, an AI model is deployed in production, where its ability to analyze and generate insights is put to the test. This is known as the inference phase, where the model recognizes patterns in external data to infer conclusions and predictions. That's the heart of the value of AI, so having the infrastructure necessary for real-time inference is ultimately what translates into actual business value and model viability.

Fortunately, inference isn't nearly as computationally intensive as training. However, unlike training, which is done intermittently, inference runs continuously in a production environment. Because of this, latency and service availability are top priorities. Also, inference workloads often span diverse storage, compute and network environments, requiring IT managers to carefully balance performance, energy efficiency and cost-effectiveness while ensuring that AI models remain secure and compliant with industry regulations.

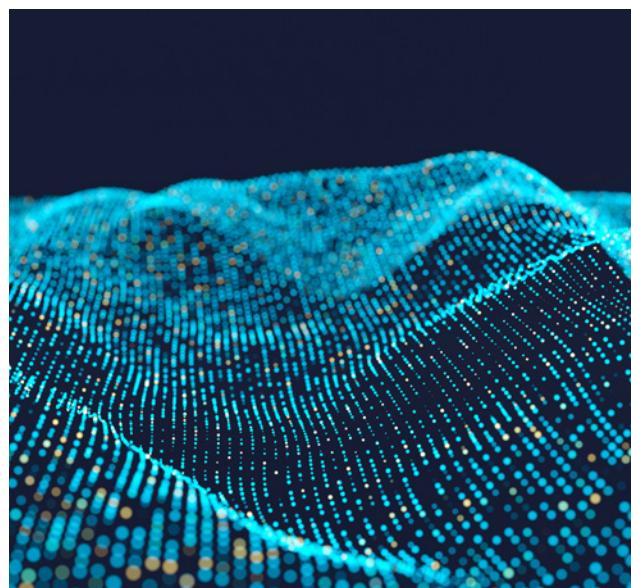
### Key challenges in AI inference deployment

- Service availability: To be useful, inference workloads must process data and return results in milliseconds, especially in mission-critical applications, such as real-time analytics.
- Resource efficiency: Inference workloads often run continuously at varying levels of demand, making dynamic scaling necessary to mitigate the risks of under- or overprovisioning.
- Data localization: Some workloads must comply with regional data laws and industry-specific regulations, thus affecting where and how AI models are deployed.

### Step 1. Optimize AI compute for instant decision-making

Unlike training, where workloads run in massively parallelized GPU environments, inference workloads are better suited AI-optimized CPUs or even edge-based processing units. That said, computational requirements vary widely, depending on the size of the model, the number of end users and the availability requirements. For instance, a model designed for fraud detection will likely need round-the-clock availability, while a customer service chatbot may need to serve thousands of real-time queries every second during peak times.

While CPUs and GPUs are both suitable for AI inference workloads, specialized low-power accelerators are often more efficient thanks to their higher dedicated memory and bandwidth. For edge inference — where AI runs locally on devices or gateways instead of remote data centers — compute solutions must also meet additional constraints, such as power efficiency, physical footprint and thermal limits, all without compromising performance.



## **Step 2. Right-size inference workloads for cost-efficient scaling**

To be useful in real-world applications, AI inference heavily depends on the rapid movement of data between compute, storage and end-user environments. An optimized network infrastructure is essential for making that possible and, in doing so, preventing latency or bandwidth bottlenecks from disrupting user experiences. For instance, SmartNICs offer software-defined hardware acceleration to ensure faster packet handling and reduced latency, while high-speed interconnects such as PCIe or CXL accelerate data movement between storage and compute infrastructures.

Unlike training, inference workloads fluctuate based on real-time demand, hence why you also need general-purpose infrastructure. As such, static provisioning can lead to wasted resources during low-usage periods or performance bottlenecks during peak usage. For greater cost efficiency and reliability, AI data centers can make use of dynamic resource allocation, workload optimization and hybrid deployment models. Right-sizing inference workloads helps achieve an optimal balance between real-time efficiency and cost control.

## **Step 3. Factor in data localization and regulatory demands**

Regulatory requirements, such as GDPR, CCPA or HIPAA, mandate that certain types of data be stored and processed within specific geographic regions. Since AI inference often involves processing potentially sensitive data in real time, enforcement of data-residency rules is a key consideration when building an AI data center. For instance, an AI application used in a healthcare organization may process patient data, which due to its sensitive nature must not be sent to an offshore cloud provider. In some cases, data localization and regulatory requirements may apply to training data as well, though this data should typically be deidentified and anonymous anyway.

Different deployment models can affect data localization and regulatory compliance in a variety of ways. For instance, if inference happens only on premises, then you have complete control over your data, but operational costs are often much higher. By contrast, a cloud data center provides scalability, but it may also raise concerns over where your data physically resides. Particularly for use cases involving sensitive data, IT leaders may instead opt for an edge AI system, where inference happens on local devices, while models and training data reside in a remote data center.



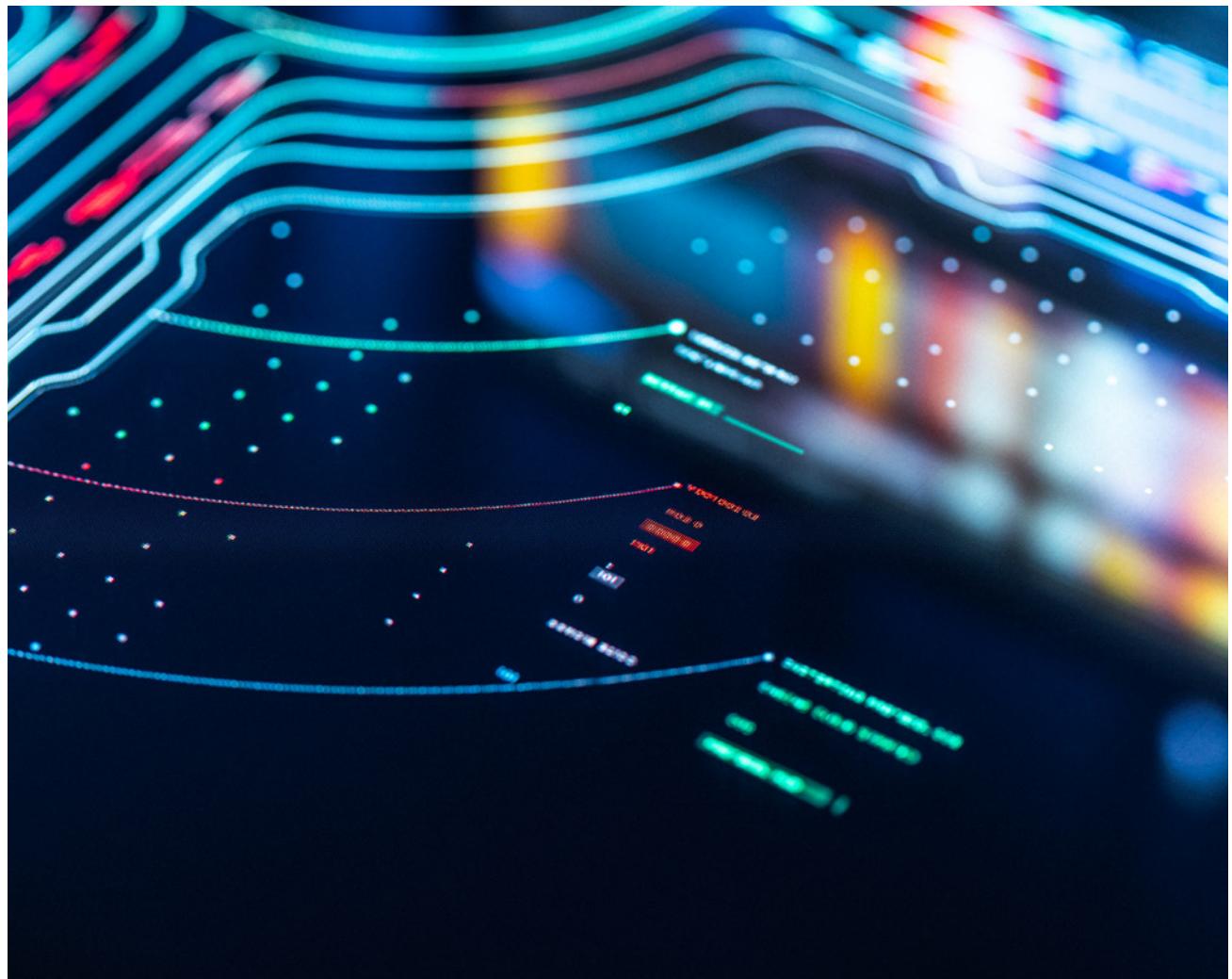
## Choosing the right AI infrastructure partner

AI is radically changing how businesses operate, compete and innovate, but success hinges on having a high-performance data center. As AI workloads continue to grow in scale and complexity, integrating specialized storage and compute and efficient networking are paramount for ensuring long-term sustainability and growth.

Working with the right AI infrastructure partner grants you access to the expertise, hardware and software

ecosystem needed to accelerate AI adoption while ensuring cost efficiency, performance and scalability. Being at the forefront of AI data center innovation, AMD delivers an end-to-end portfolio of AI solutions — from CPUs and GPUs, such as AMD EPYC™ and AMD Instinct™, to advanced networking solutions and even AI PCs — you need to build the AI-ready data center of the future.

[Learn more about AMD's cutting-edge AI solutions](#)



## References

1. The Root Causes of Failure for Artificial Intelligence Projects and How They Can Succeed: Avoiding the Anti-Patterns of AI | RAND
2. Number of Parameters in GPT-4 (Latest Data)
3. How 25,000 Computers Trained ChatGPT | Towards Data Science