White Paper

# Successfully Moving Enterprise GenAI from Proof of Concept to Production

Peter Rutten          Madhumitha Sathish
March 2025

## IDC OPINION

Generative AI (GenAI) is a transformative technology. It is being used for chatbots, coding, customer service, even R&D to generate new materials or product designs. IDC has identified hundreds of use cases for GenAI.

In the mere two years since GenAI went mainstream, hundreds of thousands of large language models (LLMs) have been developed by enterprises, hyperscalers, academia, and the open source community. These LLMs can generate anything from art to marketing copy to programming code to molecular structures to chip designs.

Many organizations have started GenAI proof-of-concept (POC) initiatives to determine if and how to invest in GenAI and generate business value. On average, organizations ran 10-20 GenAI POCs in 2024. About 50% of those had moved to production by early 2025, while 25% failed and 25% are still in the POC stage.

POCs, of course, allow an organization to quickly learn whether its use cases are feasible, whether they will provide business value, and whether they fit in their technology environment. The GenAI POC failure rate shows how challenging it can be to successfully complete such a POC.

Even more challenging is the next step. IDC now sees many businesses attempting to move from POC to production. They train (or fine-tune) a GenAI model at scale and then integrate the AI functionality into existing business processes for inferencing, at which point the model will start to field thousands or even tens of thousands of queries per day. Many businesses find out that this POC-to-production process is not a linear event from start to finish but rather a continuous cycle.

Scaling a GenAI POC to production is a complex process that requires planning, testing, and meticulous implementation. There is a phenomenon that is sometimes referred to

as "POC limbo," where the effort it takes to move from POC to full-service production is underestimated and leads to a GenAI initiative dying on the vine. When that happens, the desired business value that AI might deliver remains out of reach because the environment is not fit for operation in a real-life setting.

Yet business leaders at many organizations are anxious to have GenAI production systems in place and are short on patience to let POCs drag on for too long. They want to see value in the form of new customer experiences, new products, new processes, or improved productivity as soon as possible. Thus it is critical for all stakeholders, including these business leaders themselves, to understand what it takes to successfully move a GenAI initiative to completion. In this paper, we will go through the required steps.

## SITUATION OVERVIEW

## The Proof-of-Concept Stage

When developing a POC, businesses create a small version of what they anticipate being the AI solution they wish to deploy. They define use cases and develop and test a model to determine whether it is viable and whether the underlying assumptions for the solution are correct.

*Recommendation: Test the most desirable use cases first, the ones that are expected to deliver the most business value based on a feasible ROI calculation.*

Despite a POC's limited scale, many decisions will be made that reverberate all the way to the production stage for training and inferencing. For example, when establishing use cases, organizations will need to decide whether agentic AI is required, which is a type of AI that makes autonomous decisions. Another example would be whether the use case might benefit from retrieval-augmented generation (RAG), wherein a vectorized database of the organization's own data is queried during inferencing for greater customization.

*Recommendation: Agentic AI is complex, requiring multiple advanced models that interact with each other, so do not start a GenAI strategy with an agentic AI use case. RAG, on the other hand, is less complicated, very popular, and useful for customizing a GenAI model.*

Critically important to establish at the start of a POC are:

- Comprehensive use cases descriptions, detailing the AI functionality and how it integrates with existing business processes
- A thoroughly mapped-out value picture that the AI functionality will be delivering to the organization, using agreed metrics — whether monetary value, productivity, or innovation
- Clear, mutually agreed expectations among all stakeholders, including IT infrastructure teams, data scientists, AI engineers, line of business, and C-level participants
- A thorough road map from POC to production, detailing the steps, anticipated timelines, and expected outcomes
- A robust but flexible (because much can still change) ROI calculation based on preliminary assumptions, such as model type, model size, model deployment, skill set requirements, and infrastructure requirements

*Recommendation: Include IT, especially the IT infrastructure team, from the very beginning — as early as the stage of defining the use case — because much of your ROI success will depend on them.*

Many organizations will start their POC with prominent proprietary LLMs from major providers, since these are easily available via APIs or plug-ins. But they are also limited in terms of industry specificity. There are plenty of highly performant open source models that are tailored for a specific use case. They also tend to be smaller (thus requiring less infrastructure) and cheaper. It's most important to ensure that development tools are available and that a model is suitable for the desired use case without the need for time-consuming trials. Ultimately, the POC needs to be viewed as an iterative process, with regular feedback loops to all the stakeholders.

*Recommendation: For general-purpose use cases the large proprietary LLMs are suitable. For more industry-specific use cases, look for tailored open source LLMs.*

Data strategies, even at an early stage are, of course, hugely important. During or immediately after deciding on the strategies, organizations must determine use case definitions, what data will be needed to train (and retrain) the model, and what enterprise data will be required to inference on with RAG. Data quality is essential for accurate, high-value, and ethical model development, and thus, data preparation is critical. Be aware that data preparation can take as long as, or longer than, the actual model training. Also, note that while most organizations have decades of experience in managing structured data, many do not have the right environment for managing unstructured data such as images, text, and video.

*Recommendation: Perform an inventory of your legacy data silos. Training data will probably be distributed across the datacenter and cloud, and it needs to be consolidated (in a data lake, for example, that can keep both structured and unstructured data) before model training in production can begin.*

Data also needs to be accessible and handled in a secure, compliant fashion that prevents proprietary or private data from exposure throughout the process of data preparation — for example, during vectorization and meta tagging. Even for the purposes of a POC, these requirements should be taken very seriously, as any less disciplined work can leak into the production stage. Regulatory frameworks, such as the EU General Data Protection Regulation (GDPR), and many other domestic and international laws can be quite punitive. Be aware that, typically, a model will be adjusted multiple times by either changing the parameters or the training data sets.

*Recommendation: Automated data input into models can help make the models more dynamic and adaptable to real-world changes.*

Often, organizations develop a pilot user interface at this stage, keeping in mind that augmenting existing user interfaces within the organization can increase adoption versus building an entirely new one.

Performance of the model will be critical, even if the POC model is not yet optimized. Performance will be a factor in how large and performant the production environment will need to be.

*Recommendation: If performance is low, identify the bottlenecks in the model right away. Once the model goes into production, a lackluster performance will drag the entire process down.*

Also critical at this stage is a cost assessment of taking the model into production, for both training and inferencing. Many factors play a role here: anticipated parameter size, data volumes, and query volume. They contribute to the size, nature, location of the hosting environment, and thus associated cost. Stakeholders in the AI initiative need a solid perspective on these costs before production starts, allowing them to scale the initiative down (or up) accordingly.

To that point, early discussions about where to host the model once it goes into production will be needed. Some developers like stateless and serverless environments, which are easily scalable and typically found in the cloud. Others prefer stateful platforms, usually behind the enterprise firewall. It depends on the anticipated size of the model, the data that it is trained on, and/or references with RAG.

*Recommendation: To have the best security, the ability to optimize performance, cost containment when anticipating rapid scaling, and the ability to leverage existing compliance policies, consider the datacenter as a hosting approach.*

Ideally, the POC stage will yield a well-prepared data set, a data pipeline, a workable GenAI model, and a good understanding of all requirements for production training.

*Recommendation: Test the model for performance, accuracy, value, security, reliability, compliance, and ethical robustness before advancing to the production stage. Also, revisit the expectations among stakeholders, including use-case definitions, value to the organization, the road map, and the ROI calculation. Adjust these based on the POC results.*

Figure 1 shows a comprehensive (albeit not exhaustive) list of steps to consider for a GenAI POC. Not all will apply to every unique project, but the list enables a reasoned conversation about what to include and exclude.

**FIGURE 1**

## Checklist for the Proof-of-Concept Stage

| Proof of Concept | |
|---|:---:|
| Bring together all stakeholders (including IT). | ☑ |
| Define one or more GenAI use cases. | ☑ |
| Evaluate whether the use case requires agentic AI. | ☑ |
| Define the value of these use cases to the organization. | ☑ |
| Develop a thorough road map for the most immediate use cases. | ☑ |
| Develop a robust ROI for the most immediate use cases. | ☑ |
| Decide on proprietary LLMs or open source models. | ☑ |
| Ensure that development tools are available. | ☑ |
| Create feedback loops to stakeholders. | ☑ |
| Determine what data will be needed for model training and RAG. | ☑ |
| Assess the ability to manage unstructured data (if needed). | ☑ |
| Alleviate legacy on-premises and cloud data silos. | ☑ |
| Evaluate the data for quality. | ☑ |
| Prepare the data for model training. | ☑ |
| Ensure data security and compliance. | ☑ |
| Train the model through multiple iterations. | ☑ |
| Adjust outcomes by changing parameters and data. | ☑ |
| Test the model for performance and identify bottlenecks. | ☑ |
| Test the model for accuracy and hallucinations. | ☑ |
| Test the model for business value. | ☑ |
| Test the model for security. | ☑ |
| Test the model for reliability. | ☑ |
| Test the model for compliance with regulatory frameworks. | ☑ |
| Test the model for ethical robustness. | ☑ |
| Develop a pilot user interface. | ☑ |
| Do a cost assessment of taking the model into production (training and inferencing). | ☑ |
| Revisit expectations among stakeholders before production. | ☑ |

Source: IDC, 2025

## Result of the POC

Your organization will now have the right expectations, an ROI, a road map, a prepared data set, a data pipeline, a workable GenAI model for training, and a good understanding of all requirements for production training.

**What will change in the next stage?** In the stage of model training in production, investments will be made into the environment that will train the model at scale.

## Model Training in Production

Now that the initiative moves to large-scale model training, an organization will need to have the necessary skill sets on board: data scientists, AI developers, infrastructure and/or cloud team members and, possibly, third-party expertise for guiding the initiative to fruition.

*Recommendation: Make sure you have IT team members that understand the unique requirements of AI infrastructure, such as cluster management, parallel file systems, low-latency networks, and GPU optimization.*

If the POC was conducted using a proprietary LLM, which may have been easier at the time, the organization will now have to reconsider this approach. These large LLMs can become difficult to scale, expensive, and slow in the production stage. Organizations may be better served by converting to smaller open source models that they customize for their specific use cases through fine-tuning or, in the inferencing stage, with RAG.

Now is also the time to make decisions about where to deploy the production system for training, whether in the enterprise datacenter, in the public cloud, in the special-purpose cloud (GPU as a service), at a colocation provider environment, or at a managed services provider (SP) location. Each has its own pros and cons. Figure 2 shows the various deployment options and how they perform on a variety of important considerations. Please note the following:

- A special-purpose cloud is a tier 2 public cloud that specializes in offering GPU instances.
- A colocation datacenter is a datacenter that provides floorspace, power, cooling, internet, and other basic requirements. An organization rents this space for its own servers and storage.
- A managed services provider provides, maintains, and manages compute and storage infrastructure for an organization's applications in the provider's own datacenter or a third-party datacenter.

**FIGURE 2**

**Considerations for the Deployment Options of a Production Training Environment**

| | Datacenter | Public Cloud | Specialty Cloud | Colocation Provider | Managed Services Provider |
|---|---|---|---|---|---|
| Opex spending | ☐ | ☑ | ☑ | ☐ | ☑ |
| Data management | ☑ | ☐ | ☐ | ☑ | ☐ |
| Available AI tools | ☑ | ☑ | ☐ | ☑ | ☒ |
| Scaling capacity | ☐ | ☑ | ☐ | ☒ | ☒ |
| Integration with other workloads | ☑ | ☐ | ☐ | ☑ | ☐ |
| Infrastructure optimization | ☑ | ☒ | ☐ | ☑ | ☒ |
| Stateless model design | ☐ | ☑ | ☑ | ☐ | ☐ |
| Security and compliance | ☑ | ☐ | ☐ | ☑ | ☐ |
| Cost | ☐ | ☐ | ☑ | ☒ | ☐ |
| Geographic distribution | ☒ | ☑ | ☐ | ☒ | ☒ |
| Vendor support | ☑ | ☒ | ☒ | ☑ | ☒ |

Legend: ☑ = good, ☒ = not good, ☐ = neutral

Source: IDC, 2025

Security and compliance (as previously shown in Figure 2) refer to encryption, access controls, and compliance with privacy laws such as GDPR or HIPAA. Production systems for model training need access policies, data governance, and a security fence, as proprietary or sensitive data is often used for training the model.

Training in production is, in many ways, all about the right infrastructure, especially when an organization decides not to train in a cloud. Training at scale in a datacenter or colocation will instantly reveal infrastructure bottlenecks that may have remained hidden during the POC stage.

*Recommendation: Size the infrastructure requirements with an infrastructure vendor. These vendors have a lot of experience with making recommendations on installed infrastructure consolidation and new infrastructure needs for AI solutions.*

Infrastructure delivery and deployment speed is a major consideration. Once the infrastructure needs have been measured and defined, the selected vendor needs to be able to deliver with speed. Sometimes, the waiting time for ordered infrastructure can be as long as six months, which is valuable time when the AI solution could have been generating revenue. It is also important to ascertain whether the vendor can deliver entire racks or only separate systems. If it can deliver racks, whether these are preconfigured and tested, it will save valuable time.

Another critical topic that a vendor should help with is the existing power and cooling situation in the datacenter. Nowadays, in extreme cases, the wattage of AI infrastructure racks can reach 100kW — five to six times the usage of a legacy general-purpose rack. Not only can this power draw be prohibitive for many datacenters but also cooling such a rack often requires technologies such as liquid cooling. An organization should choose to work with a vendor that has both rack-based and direct-to-chip liquid cooling expertise.

*Recommendation: The main culprits for this power envelope are, of course, the GPUs. The latest generation is especially difficult to obtain. Not all vendors are equally able to fulfill large orders of servers or racks with the latest generation of GPUs simply because there is not enough supply. Select a vendor with proven access to GPU supplies.*

There are three other prominent infrastructure factors: network latency and bandwidth for fast transmission of large volumes of data across systems and racks, the right storage environment with multiple tiers that include AI-optimized flash storage, and a complete software stack for managing a compute cluster.

Figure 3 shows a comprehensive (but not exhaustive) list of steps to consider for taking a GenAI initiative to training in production. As previously mentioned, every project is different, but this list can facilitate a conversation about what is needed and not needed for your organization's project.

FIGURE 3

## Checklist for Model Training in the Production Stage

| Model Training in Production | |
|---|---|
| Ensure that all required skills are available (in-house or third party). | ☑ |
| Establish alignment between the teams involved. | ☑ |
| Decide on open source models versus proprietary LLMs. | ☑ |
| Establish access policies. | ☑ |
| Develop data governance policies. | ☑ |
| In deciding on where to deploy the model (i.e., datacenter, cloud, colocation provider, or managed SP), consider the following: | |
| • Open versus token based versus capex | ☑ |
| • Compute performance (GPUs) | ☑ |
| • Power requirements and cooling | ☑ |
| • Data management | ☑ |
| • Available AI tools | ☑ |
| • Scaling capacity | ☑ |
| • Compute optimization flexibility | ☑ |
| • Stateless versus stateful | ☑ |
| • Available security and compliance | ☑ |
| • Cost | ☑ |
| • Geographic distribution | ☑ |
| Determine infrastructure requirements. | ☑ |
| Decide on a build-it-yourself or turnkey solution. | ☑ |
| Train the model through multiple iterations. | ☑ |
| Adjust outcomes by changing parameters and data. | ☑ |
| Test the model for performance and identify bottlenecks. | ☑ |
| Test the model for accuracy and hallucinations. | ☑ |
| Test the model for business value. | ☑ |
| Test the model for security. | ☑ |
| Test the model for reliability. | ☑ |
| Test the model for compliance with regulatory frameworks. | ☑ |
| Test the model for ethical robustness. | ☑ |
| Identify the training performance bottlenecks. | ☑ |

Source: IDC, 2025

# The Result of Model Training

The organization now has a fully trained model, the environment to establish training and retraining as a continuous process, and the infrastructure to move to inferencing in production.

**What will change in the next stage?** At the stage of inferencing in production, the environment needs to be able to reliably manage many simultaneous queries with, potentially, large data inputs and outputs in near real time (in other words, with low latency). This process often calls for several systems in succession (especially with agentic AI) and, in many cases, across geographies. In other words, what will change is that the end user's expectation of immediacy when using the AI solution will become the main driver for infrastructure and other considerations.

# Transitioning to Inferencing in Production

In the inferencing stage too, decisions must be made regarding where to deploy the production system, with some of the same considerations for training (refer back to Figure 2). However, with inferencing, certain factors require a slightly different approach. For example, if the AI solution's end users consist of a few hundred employees in a campus location, then scaling the production environment is not a concern.

On the other hand, if there will be tens of thousands of end users per day (consumers, for example) with the potential of rapid expansion, then scaling becomes mission critical. Intuitively, this would suggest a cloud approach, but the cloud has proven to become very expensive very quickly in such scenarios. A datacenter strategy that can scale with the most common volume expansions and contractions, combined with a built-in cloud-bursting capability, will be more cost efficient.

*Recommendation: Prepare the infrastructure for slightly larger–than–expected numbers of simultaneous queries, as too many queries can choke the system (or users can get cut off). Also, anticipate larger data inputs or requested data outputs (e.g., large images, videos, or text files) than what was experimented with in the POC.*

Data volumes from queries may also require a more high-bandwidth, low-latency network than was needed in the previous stage because of the near-real-time response times expected. Plus, large numbers of concurrent queries will demand multilevel parallel compute capacity in the form of GPUs or other coprocessors, larger caches, additional memory, and dense cluster configurations. Load testing is highly recommended to determine how the environment handles various query volumes and sizes. And, just as with training, when bringing in GPUs into a dense compute environment, the same power and cooling requirements will play a role. Also, storage

will need to be performant, with large volumes of fast reads and writes tasking the storage environment. Automated monitoring of the infrastructure will help reveal performance bottlenecks that can then be remediated.

*Recommendation: Many organizations use the same infrastructure that they trained an AI model on for inferencing on that model. This is a cost-saving approach, but only under two conditions: the end-user queries are expected to be low volume, and the infrastructure is not needed for retraining or launching new AI model training exercises.*

When deploying in the datacenter or at a colocation provider, decisions must be made about a build-it-yourself or turnkey solution. Build-it-yourself solutions require specialized IT skill sets for building and configuring clusters of compute with accelerators, load balancers, schedulers, parallel file systems, and so forth. Full-stack turnkey solutions are optimized to run enterprise-scale AI workloads. They include compute, storage, and networking coupled with the necessary infrastructure software stacks and managed by an automated, unified control plane for activities such as deployment, configuration, operations, scaling, observability, security, and role-based access.

*Recommendation: Once ready to roll out the service, it is recommended to do so in stages, not all at once. Start with limited numbers of end users and with just a subset of features, monitor how the service performs, evaluate results, then start scaling by increasing end users and/or expanding feature sets. Different approaches exist as well, for example, deploying different versions of the service to different end-user groups, then deciding which one performs better.*

Production systems for inferencing also need access policies, data governance, and security — often proprietary or sensitive data is used for RAG, for example. Typically, these enterprise-grade requirements can be achieved by integrating the systems with the rest of the datacenter landscape. Integration is important not just from an IT perspective but also from a business processes perspective. Integration with the existing landscape, workloads, and businesses processes can be achieved with APIs but is sometimes complex, requiring third-party expertise.

*Recommendation: It is advisable that GenAI functionality is not standalone with its own user interface but rather gets integrated with the software that end users already use daily. This enhances end-user adoption.*

Integration with the existing landscape also promotes reliability. GenAI solutions need to run on platforms that can guarantee all the classic reliability, availability, and serviceability (RAS) features that general-purpose workloads have required for decades. This means that standard approaches such as redundancy and failover need to be applied to the AI environment to achieve the expected SLAs.

Finally, there are several best practices that will help with evolving the organization from one that runs POCs to one that operates as a mature AI-augmented business:

- If MLOps were not yet started in the POC or training stages, in this stage, it will become essential to efficiently manage and deploy the AI model as well as continuously monitor and evaluate it for improvements and updates.
- Implementing continuous integration/continuous deployment (CI/CD) provides for automated testing, code change integration, and deployment.
- AI life-cycle management will become important as many models and many versions of the same model will proliferate. This requires clear processes to manage, update, or archive them.
- Monitoring the model (not the infrastructure, in this case) for performance and accuracy in production is important as models will slowly deteriorate because of the queries they act on. Thus, sometimes, they require retraining. Model monitoring can be automated, and it is recommended to install feedback loops that enable the model to retrain itself for better accuracy.

*Recommendation: Design the entire AI POC and production process (both training and inferencing) to be repeatable so that introducing a new use case and launching a new AI model will be fast and efficient.*

Figure 4 shows a comprehensive (but not exhaustive) list of steps to consider for taking a GenAI initiative to inferencing in production.

**FIGURE 4**

## Checklist for Model Inferencing in the Production Stage

| Model Inferencing in Production | |
|---|:---:|
| Decide on the type of staged model rollout. | ☑ |
| Assess the anticipated amounts of simultaneous queries. | ☑ |
| Assess the anticipated size of data inputs and data outputs. | ☑ |
| In deciding on where to deploy the model (i.e., datacenter, cloud, colocation provider, or managed SP), consider the following: | |
| • Open versus token based versus capex | ☑ |
| • Compute performance (GPUs) | ☑ |
| • Power requirements and cooling | ☑ |
| • Data management | ☑ |
| • Available AI tools | ☑ |
| • Scaling capacity | ☑ |
| • Compute optimization flexibility | ☑ |
| • Stateless versus stateful | ☑ |
| • Available security and compliance | ☑ |
| • Cost | ☑ |
| • Geographic distribution | ☑ |
| Determine infrastructure requirements. | ☑ |
| Decide on a build-it-yourself or turnkey solution. | ☑ |
| Adjust access policies. | ☑ |
| Adjust data governance policies. | ☑ |
| Integrate the solution with the rest of your landscape. | ☑ |
| Perform load testing on query volumes and sizes. | ☑ |
| Install RAS features. | ☑ |
| Establish an MLOps practice. | ☑ |
| Establish continuous integration/continuous deployment. | ☑ |
| Establish AI life-cycle management. | ☑ |
| Establish monitoring systems for model degradation. | ☑ |
| Install feedback loops for retraining the model. | ☑ |
| Establish automated infrastructure monitoring. | ☑ |

Source: IDC, 2025

# CONSIDERING SUPERMICRO AND AMD

In the rapidly evolving landscape of AI, enterprises face significant challenges when moving AI initiatives from proof of concept to full-scale production. Supermicro, in collaboration with AMD, provides comprehensive solutions that enable businesses to overcome several hurdles (as previously discussed in this paper). By offering high-performance infrastructure, optimized cooling, and end-to-end deployment services, Supermicro and AMD ensure seamless scalability and operational efficiency for AI workloads.

Supermicro supports AI at scale with innovative servers optimized for AI training and inferencing. Systems featuring AMD EPYC processors and AMD Instinct MI325X/MI300X GPUs provide the computational power needed for deep learning models. Servers such as the air-cooled AS-8126GS-TNMR and liquid-cooled AS-4126GS-TNMR2-LCC — with two AMD EPYC 9005 CPUs and eight AMD Instinct MI325X GPUs — exemplify this approach, featuring robust compute capabilities for AI training and inferencing in demanding enterprise environments. These configurations are designed for high throughput and low latency, ensuring the efficient execution of AI workloads.

To meet enterprise AI scalability needs, Supermicro offers full-stack solutions covering compute, storage, networking, and cooling. These end-to-end solutions facilitate seamless integration of AI workloads, ensuring compatibility across on-premises, cloud, and colocation environments. Supermicro's AMD EPYC–based servers with Instinct MI325X and MI300X GPUs ensure that enterprise AI infrastructure can be efficiently deployed at scale.

The company's server lineup features the AS-8126GS-TNMR and AS-4126GS-TNMR2-LCC models, both equipped with AMD Instinct MI325X GPUs. These GPUs are optimized for demanding tasks like high-performance computing (HPC), AI applications including deep learning training, industrial automation, retail, finance, analytics, life sciences, and climate and weather modeling. These are high-performance GPU servers designed to meet the demands of AI, deep learning, LLMs, and HPC workloads. Supermicro offers a wide variety of GPU server models to accommodate options in cooling, I/O expansion, and local NVMe storage.

The AS-8126GS-TNMR and the AS-4126GS-TNMR2-LCC both support up to eight MI325X GPUs, while the AS-8125GS-TNMR2 and the AS -4125GS-TNMR2-LCC support eight MI300X accelerators and offer high memory capacity, fast data access through NVMe drives, and exceptional network bandwidth with 400GbE QSFP ports, making them powerful solutions for complex, large-scale AI workloads. The AS-8126GS-TNMR supports dual AMD EPYC 9005/9004 processors and eight AMD Instinct MI325X GPUs offering robust processing power and scalability. The AS-8126GS-TNMR and the AS-

4126GS-TNMR-LCC are powered by dual EPYC 9575F high-frequency processors with 64 cores and feature eight MI325X GPUs, providing a similar performance that is optimized for faster deployment with a 24-hour shipping turnaround.

AI infrastructure demands significant power and efficient cooling solutions. Supermicro addresses these challenges with a range of cooling solutions — including liquid-cooled server options, cold plates, CDUs, and outdoor heat exchangers — thus ensuring optimal system performance. While customers handle plumbing and piping installation, the company assembles and tests entire systems before deployment and ensures that AI datacenters can implement energy-efficient infrastructure for high-density workloads.

Advanced liquid cooling solutions, such as those seen in AS-4126GS-TNMR2-LCC, optimize thermal management, allowing enterprises to repurpose power budgets by consolidating legacy hardware and allocating resources for AI workloads. The AS-4126GS-TNMR2-LCC also features liquid cooling, providing superior thermal management for high-density configurations and making it ideal for environments that require consistent performance under heavy workloads, such as AI model training for extended periods. In addition to offering liquid-cooled GPU servers, Supermicro supplies a complete end-to-end liquid cooling solution including manifolds, pumps, management software, and the entire outdoor chiller tower.

Thus, as evident from the features of these Supermicro servers, the products are well suited for AI training and inferencing because of their powerful GPUs and ample memory, enabling fast training of complex models and efficient inferencing for real-time AI applications. The AMD Instinct MI325X and MI300X accelerators ensure that these servers can handle memory-intensive tasks such as deep learning, large language models, and AI inference with ease.

For enterprise AI, these solutions provide scalability and flexibility, allowing businesses to integrate AI into their existing infrastructure for predictive analytics, real-time decision-making, and data-driven applications. The high-performance configurations make these servers both ideal for transitioning from GenAI POC to production — supporting the computational power needed for large-scale AI model deployment and continuous inference — and essential for the development and deployment of next-gen AI applications.

## Comprehensive Testing and Implementation Support

Supermicro also offers factory-based testing, ensuring that all hardware and software components function optimally before deployment. This includes software preinstallation, network debugging, and testing at both the rack level and the multirack level. This rack-level and multirack-level integration process ensures that AI systems are

optimized before reaching customer datacenters, thus reducing deployment risks and improving time to market.

## Efficient Deployment Through Multiple Manufacturing Sites

A key advantage of Supermicro is its multiple manufacturing locations spread across San Jose (California), the Netherlands, Malaysia, and Taiwan, which enable faster and more efficient deployment of AI servers. By distributing production across different regions, the company has the potential to minimize supply chain disruptions and reduce lead times. This manufacturing location strategy enhances scalability and allows enterprises to deploy AI infrastructure seamlessly across global datacenters, colocation facilities, and on-premises environments. It also ensures that enterprises have steady AI operations, even in the face of logistical challenges.

## AI Solutions Tailored for Industry-Specific Needs

Supermicro serves all verticals, including some of the most demanding industries such as financial services, retail, and telecommunications. The company's vertical-specific expertise ensures that enterprises in all sectors can leverage AI infrastructure that is tailored to their performance, security, and compliance needs.

By offering an extensive portfolio of AI-optimized servers, robust testing protocols, and scalable solutions, by continuously updating product offerings, and thus by ensuring that organizations have access to the latest AI hardware innovations, Supermicro addresses several needs of enterprises that transition AI from POC to production. The solutions cater to the critical importance of agility, speed, time to market, and beyond, to the concept of "time to day one" where systems are operational on the first day of deployment. Scalability is another vital factor, as AI applications demand infrastructure that can adapt quickly to changing demands. To address these needs, Supermicro offers solutions that can be produced at rack scale, providing customers with the flexibility to rapidly adjust their infrastructure. This capability allows businesses to swiftly shift between different configurations, responding to new opportunities in AI as they arise. Thus Supermicro, in collaboration with AMD, is well positioned as an agile infrastructure partner.

## CONCLUSION

GenAI has rapidly evolved, with untold numbers of LLMs developed for a multitude of applications. Many organizations have initiated GenAI POC projects, but transitioning from POC to production remains challenging. Successful scaling requires meticulous planning, testing, and implementation. Key steps include defining comprehensive use cases, establishing clear expectations, and involving IT infrastructure teams early. Data

quality and preparation are critical, and organizations must consider infrastructure needs such as power and cooling requirements.

Supermicro and AMD offer solutions to address these challenges, providing high-performance servers optimized for AI training and inferencing. Their systems ensure efficient execution of AI workloads. Supermicro's comprehensive solutions cover compute, storage, networking, and cooling, facilitating seamless integration across various environments. Its multiple manufacturing locations enhance deployment efficiency.

For enterprises, these solutions offer scalability, flexibility, and the computational power needed for large-scale AI model deployment. Supermicro's rigorous testing protocols and industry-specific expertise ensure that AI infrastructure meets performance, security, and compliance needs. By addressing the critical importance of agility and speed, Supermicro and AMD position themselves as key partners for organizations transitioning AI from POC to production.

## ABOUT IDC

International Data Corporation (IDC) is the premier global provider of market intelligence, advisory services, and events for the information technology, telecommunications, and consumer technology markets. With more than 1,300 analysts worldwide, IDC offers global, regional, and local expertise on technology, IT benchmarking and sourcing, and industry opportunities and trends in over 110 countries. IDC's analysis and insight helps IT professionals, business executives, and the investment community to make fact-based technology decisions and to achieve their key business objectives. Founded in 1964, IDC is a wholly owned subsidiary of International Data Group (IDG, Inc.).

## Global Headquarters

140 Kendrick Street
Building B
Needham, MA 02494
USA
508.872.8200
Twitter: @IDC
blogs.idc.com
www.idc.com