# FlexPod AI for retrieval-augmented generation

## FlexPod AI for RAG: save time, reduce risk, and increase security

In today's competitive landscape, businesses are increasingly turning to AI to gain insights, to automate processes, and to create new user experiences. However, deployment of AI, and particularly generative AI, can be more challenging than expected. These systems demand robust and agile IT infrastructure to handle complex data processing, to maintain security, and to comply with ever-changing regulations. The stakes are high, because inefficiencies can lead to data bottlenecks, security breaches, and compliance violations, all of which can have significant business implications such as lost revenue, diminished customer trust, and legal penalties.

Understanding these challenges, Cisco and NetApp have joined forces to develop FlexPod® AI, a converged infrastructure

solution that is specifically engineered to support the unique requirements of AI-driven applications. FlexPod AI is designed to overcome the critical pain points that businesses face when they implement AI solutions. FlexPod AI helps organizations meet these challenges by integrating retrieval-augmented generation (RAG) technology, which significantly enhances the quality and relevance of generative AI outputs. This integration enables FlexPod AI to deliver more accurate and contextually relevant results by efficiently retrieving necessary data from a vast corpus of information.

Moreover, the FlexPod AI architecture is built with a focus on streamlining compliance and bolstering security. The system's advanced data management capabilities confirm that regulatory requirements are met without sacrificing performance. Businesses can embrace AI confidently, knowing that their infrastructure is not only

powerful and efficient, but also compliant with the latest industry standards.

In essence, FlexPod AI for RAG solutions is more than just a technological innovation; it's a strategic business enabler. It offers companies a way to overcome the inherent complexities of generative AI, with a secure, compliant, and high-performing AI infrastructure that can meet the demands of modern AI applications and deliver tangible business value.

## What is retrieval augmented generation?

RAG is an innovative approach in natural language processing (NLP) that combines retrieval-based and generation-based methods. It uses a vast corpus of data to retrieve relevant information, which is then used to generate more accurate and contextually appropriate responses. This hybrid model improves the performance of AI applications by generating more accurate content and by

increasing the relevance of that content.

Organizations should pay particular attention to the advantages of RAG, because they directly tackle some of the most pressing concerns in AI-driven operations—accuracy, reliability, and efficiency. In an era where data reigns supreme, the ability to swiftly access and use the correct information is paramount. RAG empowers systems with a more nuanced understanding of user queries, leading to outputs that are not only more relevant, but also of a higher quality. These benefits are crucial for businesses that rely on AI to interact with customers, to analyze data, or to generate content, because they can significantly enhance user satisfaction and trust.
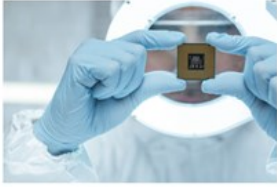
Furthermore, by reducing errors and by improving decision-making processes, RAG can help companies avoid costly mistakes and make better-informed strategic choices. The value proposition of RAG lies in its potential to transform vast, unstructured data landscapes into actionable insights, promoting innovation, a competitive edge, and bottom-line results for businesses that take advantage of its capabilities.

By enabling smarter decision-making and by automating complex processes, AI technologies are transforming industries. However, the increasing complexity and volume of data require advanced solutions to manage AI workloads effectively. FlexPod for AI use cases provides a scalable, high-performance infrastructure that meets these demands. And by integrating RAG with FlexPod AI, organizations can enhance the efficiency and accuracy of their AI models. The table below is just a small sample of the many enterprise use cases for RAG on FlexPod for AI .

| Enterprise use cases for RAG | |
|---|---|
| Customer support | Automate responses to customer inquiries by retrieving relevant information from knowledge base articles. |
| Educational tools | Assist students by providing detailed explanations to answer their queries, using verified educational resources. |
| Writing assistance | Help writers by retrieving relevant information about a topic and by generating drafts or enhancing existing content. |
| Summarization | Generate concise summaries of long documents, reports, or articles by retrieving key points and synthesizing them. |
| E-commerce | Provide personalized product recommendations by retrieving user-specific data and generating tailored suggestions. |
| Entertainment | Suggest movies, books, or music based on user preferences and by retrieving relevant content from databases. |
| Scientific research | Assist researchers by retrieving relevant studies and generating summaries that support their research efforts. |
| Medical diagnosis | Support healthcare professionals by retrieving patient data and generating possible diagnoses based on similar cases. |
| Document analysis | Assist legal professionals by retrieving relevant case laws or regulations and generating summaries. |
| Compliance checks | Automate compliance checks by retrieving relevant policies and generating compliance reports. |
| Multilingual support | Provide accurate translations by retrieving contextually relevant examples and generating translations that maintain the original meaning. |
| Localization | Generate localized content that aligns with cultural and regional nuances by retrieving and adapting information appropriately. |

Accuracy


Personalization


Scalability


Multi-Lingual


Reliability

In an enterprise environment, the benefits of RAG are multifaceted and impactful. By integrating RAG into their systems, enterprises can expect a marked improvement in the quality of AI-generated content, leading to more accurate and insightful interactions with customers and data. This accuracy translates into enhanced trust and satisfaction from end users who receive reliable and contextually relevant information. Moreover, the reduction in AI "hallucinations" promotes business decisions that are informed by precise and verifiable data, mitigating risks and reinforcing the integrity of automated processes. Ultimately, RAG serves as a powerful tool in an enterprise's arsenal by streamlining operations; by supporting compliance with regulatory standards; and by providing a competitive edge through smarter, more efficient AI capabilities.

## FlexPod AI: A robust foundation for AI workloads

FlexPod AI is designed to support the most demanding AI and machine learning (ML) workloads. Its architecture combines Cisco UCS Servers, Cisco Nexus switches, and a NetApp® data management plane to deliver a unified, scalable, and high-performance infrastructure.

The integration of RAG with FlexPod AI represents a significant advancement in AI infrastructure, offering enhanced performance, scalability, and efficiency for complex AI workloads. By adopting this powerful combination, organizations can unlock new possibilities in AI applications, driving innovation and achieving a competitive advantage. The FlexPod architecture aligns with the meticulous design that is outlined in the FlexPod Datacenter with Red Hat OCP Bare Metal Manual Configuration with Cisco UCS X-Series Direct.

### Key FlexPod AI features

**Scalability:** Easy-to-scale compute and storage resources to meet growing AI demands

**Performance:** High throughput and low latency to accelerate inferencing

**Simplicity:** Simplified deployment and management with validated designs and automation tools
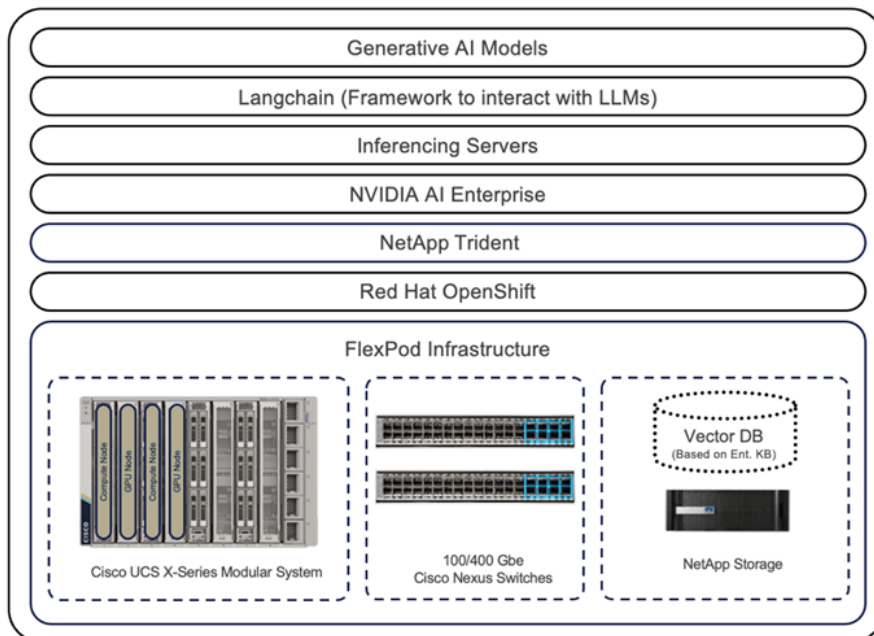
**Flexibility:** Support for various AI frameworks and tools, promoting compatibility with diverse AI applications

This foundation is further strengthened by incorporating NVIDIA AI Enterprise and NVIDIA NIM microservices, creating an optimized environment to run RAG workloads.
The platform's configuration adheres to industry best practices and prioritizes security, delivering a high-performance solution for generative AI tasks. The versatility of FlexPod AI extends beyond just generative models; it is scalable to support a wide array of AI applications, including training, fine-tuning,

and various inferencing scenarios. The key to unlocking this potential lies in appropriately sizing the platform to accommodate the specific demands of these diverse AI use cases.

## Unlocking the full potential of AI

In summary, FlexPod Datacenter for AI is a comprehensive solution that is specifically designed for AI/ML environments. By seamlessly integrating GPU, compute, storage, and networking technologies into a unified platform, it empowers AI applications while it enhances operational efficiency and agility. Notably, this solution excels in reducing complexity, offering organizations a powerful and streamlined infrastructure to unlock the full potential of AI within their operations. To find out more about RAG on FlexPod AI, review the Cisco Validated Design FlexPod for Accelerated RAG Pipeline with NVIDIA NIM and Cisco Webex.