



Docker, Putting the AI in Containers

How Containerization Can Accelerate
Enterprise Artificial Intelligence Development

By Mark Hinkle

The Strategic Imperative of AI in 2024

The winds of change are sweeping across industries, propelled by the transformative power of Generative Artificial Intelligence (GenAI). In 2024, AI has become a vital strategic imperative for enterprises seeking to stay ahead of the curve. While some may view AI with hesitation, the reality is that ignoring its potential risks them falling behind. This whitepaper aims to equip you with the knowledge and tools to unlock the transformative potential of AI, starting with the powerful platform of Docker containerization.

The Cambrian Explosion of Artificial Intelligence

We are familiar with chatbots for desktop users like ChatGPT and Google Gemini. However, the landscape of enterprise applications is teeming with examples of AI driving differentiation and success. Consider healthcare, where AI algorithms aid in early disease detection and personalized treatment plans, or finance, where AI-powered fraud detection systems and algorithmic trading are reshaping the industry. AI-driven robots optimize production lines in manufacturing, and predictive maintenance minimizes downtime. We are seeing an even more significant expansion where new types of AI systems provide solutions to problems previously not attainable with machine learning. Now, new GenAI systems provide capabilities to solve enterprise's most pressing issues faster and more efficiently than ever.

In 2023, IBM reported that 42% of IT professionals at large organizations report that they have actively deployed AI while an additional 40% are actively exploring using the technology¹. Across the board, businesses are leveraging AI to innovate, gain market share, and secure a competitive edge.

Consulting firms are now using AI to build a reputation for solving problems, which is a key part of focusing on their customers' needs."Over the next three years, Capgemini plans to invest \$2 billion in generative AI, while PwC and Cognizant will each invest \$1 billion. Companies with generative AI are partnering with hyperscalers to access computing power: OpenAI with Microsoft, Anthropic and Stability AI with Amazon, and Cohere with Oracle².

The landscape of AI models has undergone a fascinating shift in a very short time. We witnessed the initial explosion of behemoths like Amazon's GPT-3, boasting billions of parameters and impressive capabilities. These large language models (LLMs) captivated the world with their ability to generate human-quality text, translate languages, and answer complex questions.

However, the allure of size soon began to encounter limitations. The sheer scale of these models presented challenges in terms of computational resources, training costs, and environmental impact. As sustainability concerns intensified and accessibility became a priority, a new breed of AI models emerged: the small and robust models.

These models, exemplified by projects like Mixtral, Microsoft's Phi, Google's Gemini, and others, operate with significantly fewer parameters, often in the millions or even tens of millions. This reduction in size doesn't equate to a decrease in capability. These models leverage innovative architectures and training techniques to achieve: impressive performance metrics, sometimes rivaling their larger counterparts.



Not only have models increased, but there has been a growth of open source ethos in artificial intelligence. [Hugging Face](#), a repository for open source artificial intelligence software, datasets, and development tools, has seen its list of models grow to over 500,000 models of all shapes and sizes suited for various applications³. These models are ideally suited for deployment in containers that can be developed locally or in the data center.

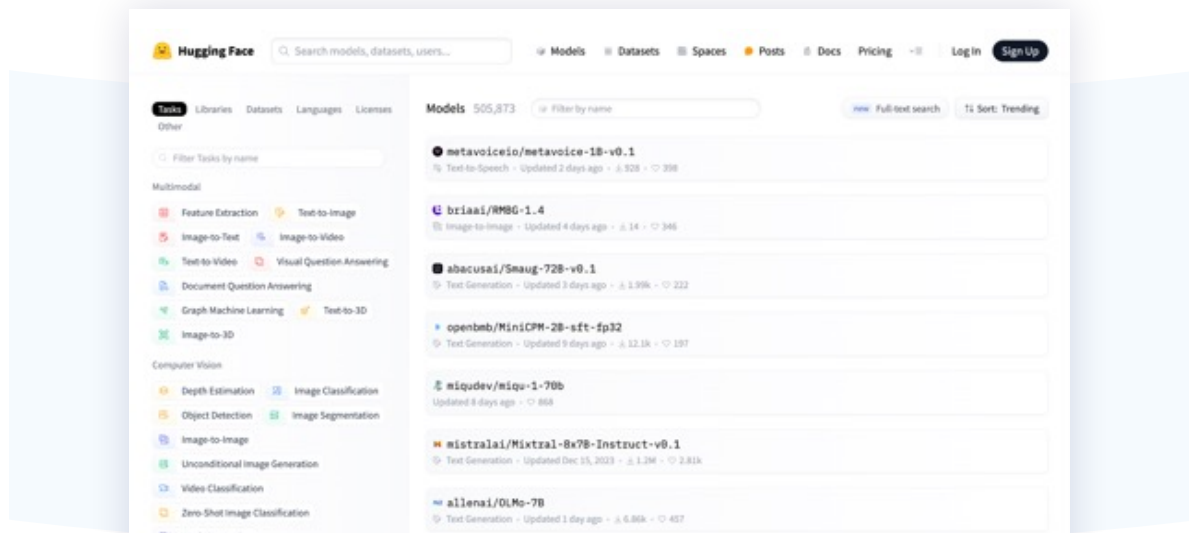


Figure 1: HuggingFace provides a repository of open source models and tools to help test and develop large language models.

This shift towards smaller, more efficient models signifies a crucial change in focus. The emphasis is no longer solely on raw power but also practicality, resourcefulness, and accessibility. These models democratize AI by lowering the barrier to entry for researchers, enterprise software developers, and even small and medium businesses with limited resources. They pave the way for deployment on edge devices, fostering advancements in areas like AI at the edge and ubiquitous computing.

They will also provide the foundation for enterprises to adapt and fine-tune these models for their usage. They'll do so using their existing practices of containerization, and they will need the tools that can provide the ability to move quickly through each phase of the software development lifecycle. Containerized software development supported by Docker is the ideal way to do this as the industry's de facto development and deployment environment for enterprise applications.

The arrival of these small and powerful models signals a new era in AI development. It's a testament to the ingenuity of researchers and a shift towards responsible and sustainable AI advancement. While large models will likely continue to play a vital role, the future of AI appears increasingly bright and diverse, fueled by the innovative spirit of these smaller, more impactful models.



Operational Drivers:

Beyond the competitive landscape, AI presents a compelling value proposition through its operational benefits. Imagine automating repetitive tasks, extracting actionable insights from massive datasets, and delivering personalized experiences with unparalleled accuracy. AI facilitates smarter, data-driven decision-making as users push projects to completion, improving efficiency, cost reduction, and resource optimization.

Alignment with Business Goals

However, more than simply deploying AI as a technology standalone, users must align AI initiatives with specific business goals and objectives. Whether driving revenue growth, expanding market share, or enhancing operational excellence, AI-driven projects can be a powerful tool when directed toward strategic priorities. For instance, AI-powered recommendation engines can boost sales, while chatbots can improve customer service, ultimately contributing to overall business success.

Digital Transformation

Moreover, AI is a cornerstone of digital transformation initiatives. Businesses are undergoing a fundamental shift towards data-driven, interconnected operations, and AI plays a critical role in unlocking new opportunities and accelerating this transformation. From personalized marketing campaigns to hyper-efficient supply chains, AI empowers organizations to adapt to ever-changing market dynamics and achieve sustainable growth.

The Clear AI Imperative

In 2024, AI is not a luxury but a strategic necessity. As competitors leverage AI to innovate and gain an edge, enterprises that fail to embrace this technology risk falling behind.

Containers: Fueling the Evolution of AI Development

Software development has faced its share of challenges. Inconsistent environments, cumbersome deployments, and complex configurations impeded progress and stifled innovation. But then came containerization, propelled by Docker, and the landscape shifted dramatically.

The development standard set by containerization will benefit the development of new artificial intelligence systems.



From Bottlenecks to Breakthroughs

Before containers, building and deploying AI applications felt like assembling a puzzle with missing pieces. Choosing and testing embeddings with multiple vector databases, bring up machine learning libraries like Tensor Flow and Pytorch, and then deploying LLMS are all part of the GenAI ecosystem but without a tool that allows you do so locally and then deploy to testing and production makes an already daunting task untenable. Docker's commitment to providing a GenAI Stack changed the game. Enabling the packaging of entire applications (including their libraries and dependencies) into portable units called containers ushered in a new era of:

- **Effortless Reproducibility:** Gone are the days of inconsistent environments leading to unpredictable results. Containers guarantee identical conditions across development, testing, and production, ensuring reliable, reproducible AI workflows.
- **Streamlined Collaboration:** Using containers, enterprises can save wasted time troubleshooting compatibility issues. Containerized AI applications offer seamless collaboration, empowering teams to build and iterate swiftly with confidence.
- **Agile Experimentation:** In lieu of tedious setup processes, Docker lets developers spin up and scale AI environments in seconds, paving the way for rapid experimentation. Devs can focus on exploring and selecting different frameworks and tools that foster project innovation.
- **Seamless Deployment:** Leverage containers to eliminate infrastructure challenges. Containerized AI applications seamlessly deploy across diverse platforms and cloud environments, minimizing friction and speeding up time-to-market.
- **Resource Optimization:** Maximize your hardware potential. Containers share the host operating system, leading to efficient resource utilization and reduced infrastructure costs.

Unlocking the Potential of AI with Containers

This shift translates into real-world benefits for AI development. This applies not only to developer onboarding which kickstarts the development process but also provides value throughout the development cycle. Moving to adding GenAI to the mix will only create another development pipeline that users need to address. Docker provides those benefits and many more:

- **Easily Deployed Development Environments:** One big advantage to using Docker whether for AI applications or others is the ability to bring up complex development environments quickly then deploy software throughout the development, testing, and deployment lifecycle.
- **Faster Development Cycles:** Streamlined workflows and reduced configuration complexities accelerate development cycles, bringing AI ideas to fruition quicker.
- **Enhanced Collaboration:** Shared, reproducible environments create a platform for smooth collaboration, boosting team productivity and innovation.



- **Dynamic Scalability:** Easily scale containerized AI applications to handle fluctuating workloads, ensuring seamless performance under pressure.
- **Cost-Effective Operations:** Efficient resource utilization minimizes infrastructure costs, allowing organizations to achieve more with less.
- **Experimentation Without Boundaries:** Explore different AI frameworks and tools without impacting existing environments, opening doors to new possibilities.

Beyond mere tools, containers have become fundamental to the evolution of AI development. They empower organizations to overcome long-standing challenges, accelerate innovation, and unlock the true potential of AI. In the following sections, we'll delve deeper into specific trends and applications that showcase the powerful synergy between containers and AI, paving the way for a future where AI is more accessible, scalable, and impactful.

Key Trends Shaping AI and Containerization

The synergy between AI and containerization is fueled by exciting, evolving, and intertwining trends:

1. **Generative AI:** This rapidly evolving field, capable of creating text, images, and even code, holds immense potential for automating tasks and fostering creativity. Docker plays a crucial role in managing the complex pipelines and diverse dependencies required for training and deploying generative AI models.
2. **Ethical AI Frameworks:** Ethical considerations emerge as AI applications become more powerful. Containerization allows organizations to implement modular, transparent, and auditable AI processes, ensuring responsible development and deployment.
3. **AI at the Edge:** Processing data closer to its source offers lower latency and improved privacy. Containerized AI applications are ideally suited for edge environments due to their portability, resource efficiency, and offline capabilities.
4. **Cloud-Native AI Development:** Migrating AI workloads to the cloud offers scalability and flexibility. Containerization facilitates seamless integration with cloud-native platforms, simplifying deployment and management across hybrid and multi-cloud environments.
5. **MLOps and Continuous Delivery:** Streamlining the machine learning (ML) lifecycle is crucial for efficient AI development. Containerization integrates seamlessly with MLOps practices, enabling continuous integration, continuous delivery (CI/CD), and automated workflows for AI pipelines.



6. **Security Considerations:** Securing AI models and protecting sensitive data are paramount. Containerization enhances security by providing isolation and granular control over access, resource usage, and network communication.
7. **Democratization of AI:** Containerization lowers the barrier to entry for AI experimentation and development by providing accessible, reproducible environments. This democratization allows diverse teams and individuals to contribute to AI innovation.
8. **Open Source Collaboration:** The open source nature of both AI and containerization foster rapid development and innovation. This collaboration benefits the entire community, creating a vibrant ecosystem of tools, libraries, and best practices.

Organizations can position themselves at the forefront of AI development by understanding these trends, the opportunity they pave for enhanced business impact, and their connections to containerization. Leveraging Docker's capabilities effectively empowers them to overcome challenges, unlock new possibilities, and drive their AI initiatives forward.

Docker GenAI Stack: Simplifying AI Integration for Modern Developers

While containers addressed the challenges associated with application component management and AI-driven tools revolutionizing code generation, a void still existed. How does one ensure the AI-driven code generation tools know the environments that Docker containers encapsulate? How can these tools be optimized to generate code that's not just functional but also optimized for containerized environments?

Docker makes it possible to bridge this gap. By infusing AI capabilities directly into the Docker ecosystem, Docker equips developers with tools that enhance both application development and containerization.

Building software with AI isn't futuristic anymore; it's the new normal. But juggling different tools and complexities can be a headache. That's where the Docker GenAI Stack comes in, a key development tool built with industry leaders.

Think of it as your one-stop shop for seamless AI integration. This unified platform, tailored specifically for generative AI, is Docker's answer to the growing demand for smooth, hassle-free AI development. No more cobbling together different pieces – the GenAI Stack streamlines everything, setting a new bar for building AI-powered applications.



Docker GenAI Stack: Key Features and Capabilities

The Docker GenAI Stack empowers developers to integrate AI seamlessly into their applications by providing comprehensive tools. It streamlines the entire process, making complex tasks manageable and intuitive. Forget about juggling disparate technologies. Let the GenAI Stack guide you step-by-step towards crafting intelligent applications.

Quick Start

With the GenAI Stack, Docker ensures developers can initiate their AI projects promptly, eliminating the typical roadblocks associated with blending diverse technologies.

Comprehensive Component Suite

The stack comes bundled with a range of pre-configured components primed for the immediate development of artificial intelligence applications. This suite includes:

Large Language Models (LLMs): A selection of models from Ollama like Llama 2, Code Llama, and Mistral, and widely recognized models such as GPT-3.5 and GPT-4.

Advanced Databases: The Docker GenAI Stack isn't just streamlining AI integration – it's turbocharging model precision thanks to powerful Neo4j vector and graph databases built right in. Forget limitations. Neo4j empowers your AI and ML models to think sharper with fewer hallucinations and deliver results you can trust.

LangChain Framework: A critical tool in the stack, LangChain ensures a smooth interface between LLMs, applications, and databases, facilitating seamless interactions. LangChain is a framework that allows developers to create agents capable of reasoning about issues and breaking them down into smaller sub-tasks. By building intermediary stages and chaining complex commands together, you can add context and memory to completions using LangChain. LangChain plays a vital role in the technology stack by enabling seamless communication between language models (LLMs), applications, and databases. It is a framework that empowers developers to establish intelligent conversational agents capable of understanding and breaking down complex problems into smaller, manageable subtasks. Through the creation of intermediate stages and the chaining together of intricate commands, LangChain provides the ability to infuse completions with context and memory, enhancing the overall quality and relevance of the generated responses.

"Everything changed this year, as AI went from being a field for specialists to something that many of us use every day. The tooling landscape is, however, really fragmented, and great packaging is going to be needed before general broad-based adoption by developers for building AI-driven apps really takes off. The GenAI Stack that Docker, Neo4j, LangChain, and Ollama are collaborating to offer provides the kind of consistent unified experience that makes developers productive with new tools and methods so that we'll see mainstream developers not just using AI, but also building new apps with it."

James Governor

Principal Analyst and Co-Founder of RedMonk



Components of the Docker GenAI Stack

The Docker Generative AI Stack, or Docker GenAI Stack, is a comprehensive suite of tools and resources designed to facilitate the development and deployment of Generative AI applications.

Let's dive deeper into its components:

Large Language Model(LLM) Deployment with Ollama

Docker's suite includes a diverse range of pre-configured open source LLMs tailored to meet the varying needs of applications. At the heart of the system, these LLMs function as the primary engines, generating text that mirrors human-like articulation, contingent on the input they process.

Llama 2 (technically not open source but with many of the same benefits), Code Llama, and Mistral stand out due to their capabilities and widespread adoption. Docker facilitates access to elite private models for developers aiming for cutting-edge solutions, notably OpenAI's GPT-3.5 and GPT-4.

Recognizing developers' challenges with LLMs, Docker introduced Ollama as part of the stack. Ollama is a dedicated Docker image that enables developers to deploy large language models quickly and easily.

Open Source LLMs hold immense potential, but getting started can be a technical maze. Don't let complex setups and troubleshooting hinder your exploration! Ollama streamlines the entire process of setting up LLMs locally. Gone are the days of wasted time and frustrating hiccups. Ollama expedites the setup, ensuring a smooth transition into the exciting world of GenAI.

Ollama doesn't just get you started. It equips you with the right tools and clear guidance to confidently navigate the LLM landscape. Whether you're a seasoned developer or a curious newcomer, Ollama empowers you to unlock the creative possibilities of Generative AI.

Neo4j Database

The Docker GenAI Stack incorporates Neo4j, the leading graph database. Graph databases, such as Neo4j, offer distinct advantages in AI development, primarily due to their ability to model, store, and query relationships between data points efficiently. This capability is especially beneficial in AI and machine learning (ML) projects that involve complex relationships and interconnected data.

Neo4j's dual capability allows for both graph-based and native vector searches, catering to diverse data retrieval needs. One of Neo4j's standout features is its proficiency in unearthing explicit and implicit data patterns and relationships. This depth of insight is invaluable for AI-driven applications. Neo4j is pivotal in data retention beyond just data retrieval. It acts as a reservoir, holding onto vital insights and ensuring that AI/ML models can recall and build upon past learnings.



Neo4J, a graph database, offers a structured representation of information, making it beneficial when integrated with a Large Language Model (LLM). This structure aids the LLM in organizing knowledge logically and coherently, with clearly defined relationships and entities that enhance response accuracy and context relevance. Knowledge graphs enhance an LLM's semantic understanding by elucidating relationships and hierarchies between concepts, improving generalization across tasks and domains. Furthermore, they provide context awareness, allowing the LLM to understand the nuances of a query better, thus enhancing the relevancy and precision of its responses.

Large Language Models adapt and learn continuously. With their structured layout, knowledge graphs offer a foundational framework for accumulating and organizing new data over time, with the flexibility to seamlessly incorporate new nodes and relationships. Neo4J's ability to integrate data from diverse sources presents a unified view, enriching the LLM's understanding of information. The graph-based nature of Neo4J also facilitates sophisticated querying and reasoning, enabling the LLM to address complex queries and execute advanced reasoning tasks. Additionally, functions like recommendation systems and anomaly detection benefit from the graph's structured nature, and the traceability and explainability provided by knowledge graphs are vital for real-world applications, ensuring that the LLM's responses are transparent and accountable.

Langchain

At the intersection of the LLM, application, and database lies LangChain. It acts as the linchpin, ensuring seamless communication and interaction among these components.

One of LangChain's primary responsibilities is overseeing vector indexing. This process is crucial for swift and precise data retrieval, especially when dealing with vast datasets typical in AI applications.

Beyond just data management, LangChain offers a robust framework. This structure empowers developers to create applications that are not only context-aware but can also reason in line with the insights provided by LLMs. Incorporating this type of AI middleware ensures that the applications are data-driven and insight-driven.

Supporting Tools and Resources

Docker's philosophy extends beyond providing technological solutions. It recognizes that developers need comprehensive tools and resources for their AI development projects to flourish.

Code Templates

Docker offers a suite of [ready-made code templates](#) (Support Bot, Stack Overflow Loader, PDF Reader and others). These are not just generic structures but are tailored for AI development, ensuring that developers can hit the ground running without getting bogged down by foundational coding.

How-to Guides

Navigating the world of AI can be daunting, especially for those new to the domain. Docker's how-to guides act as a compass, providing clear, step-by-step directions for various processes, ensuring that developers always have a reference point.



Best Practices

Docker curates a list of best practices for Generative AI development. These guidelines ensure that developers create efficient applications and adhere to industry standards, ensuring longevity and relevance.

The GenAI Stack actively orchestrates components like LangChain and various supporting tools, creating a nurturing environment for developers. Imagine having everything you need right at your fingertips: cutting-edge tech, best practices, and a community pushing the boundaries of what's possible. That's the GenAI Stack difference.

The seamless integration of LLMs, Ollama's support, and the Neo4j database form the backbone of Docker's AI-driven development approach. LLMs generate the content, Ollama ensures smooth deployment and integration, and Neo4j provides the data foundation. This trinity ensures that developers have a holistic, efficient, and powerful toolset, enabling them to harness the full potential of Generative AI in their applications.

The Docker GenAI Stack is a holistic solution for developers venturing into the realm of Generative AI. It provides the necessary technological tools, support, and resources to ensure successful project execution.

"We're all here to help teams close the gap between the magical user experience GenAI enables and the work it requires to actually get there. [The Docker GenAI Stack] is a fantastic step in that direction."

Harrison Chase
Co-Founder and CEO of LangChain

Components of Docker for AI: GPU Access for AI Applications

AI applications often require significant computational resources like GPUs to train and run models. Docker's ability to access GPUs allows developers to create and deploy AI applications consistently and reproducibly, regardless of the underlying hardware. Direct hardware access means that developers can create and test AI applications on their local machines and then deploy them to production environments with confidence that they will work as expected.

To use Docker for AI development, developers can create a Docker image that includes all the necessary dependencies, such as libraries like TensorFlow or PyTorch, and then use that image to create a container that can access the GPU. The resulting container can then be utilized to train and run AI models.

Docker's ability to access GPUs makes it an ideal development environment for AI applications. By creating GPU-accelerated Docker containers, developers can create and deploy AI applications consistently and reproducibly, regardless of the underlying hardware. Hardware access via Docker makes testing, debugging, and scaling AI applications easier and ensures they work as expected in production environments.



Docker & AI: Wisdom of Crowds Comes to Artificial Intelligence

As AI-driven development gains traction, there's a pressing need for tools that seamlessly integrate AI capabilities into modern software development processes. How can Docker, with its containerization prowess, address this evolving demand?

The synergy between Docker and AI/ML signifies more than just technological advancement; it represents a paradigm shift in how we perceive software development. With Docker providing a consistent environment for application components and AI-driven tools, developers are poised to navigate the complexities of modern software development with unprecedented ease.

Docker GenAI Stack simplifies AI integration for developers. Equipped with various features and components, have been designed to bridge the gap between traditional software development and AI-driven methodologies.

The case underscores the transformative potential of combining Docker's containerization with AI-driven tools, as it promises to redefine software development paradigms, offering enhanced efficiency, security, and innovation. The question remains: How can organizations best leverage this synergy to stay ahead in the competitive software development arena?

"The driving force uniting our collective efforts was the shared mission to empower developers, making it very easy for them to build GenAI-backed applications and add GenAI features to existing applications."

Emil Eifrem

Co-Founder and CEO of Neo4j



The answer is Docker, the AI developer's toolkit.

The synergy between GenAI and agile containerized development promises even more innovations. Imagine AI tools that can predict application scaling requirements based on user behavior and automatically configure Docker containers to handle the load. Or consider Docker environments that adapt in real-time based on AI-driven analytics, ensuring optimal resource allocation.

The confluence of Docker and AI/ML is not just reshaping software development; it's setting the stage for a future where applications are more innovative, development is faster, and the boundaries between code, environment, and intelligence are seamlessly blurred.

The Docker GenAI Stack is available in Early Access now and is accessible from the Docker Desktop Learning Center or [on GitHub](#).

1 IBM Global AI Adoption Index 2023, "IBM", <https://newsroom.ibm.com/2024-01-10-Data-Suggests-Growth-in-Enterprise-Adoption-of-AI-is-Due-to-Widespread-Deployment-by-Early-Adopters>, retrieved February 12, 2024

2 Global Tech Spend Will Grow 5.3% In 2024, Forrester, Michael O'Grady, Forecast Analyst and Michael Kearney, Data Researcher, Jan 16, 2024, retrieved January 12

3 Hugging Face Models, <https://huggingface.co/models>, retrieved on Feb 2, 2024

