



Adria Forum 2025

On premise AI with Cloud Foundation 9 L100

Marko Žuvanić Territory Sales Bruno Šunjić System Engineer

14.10.2025.



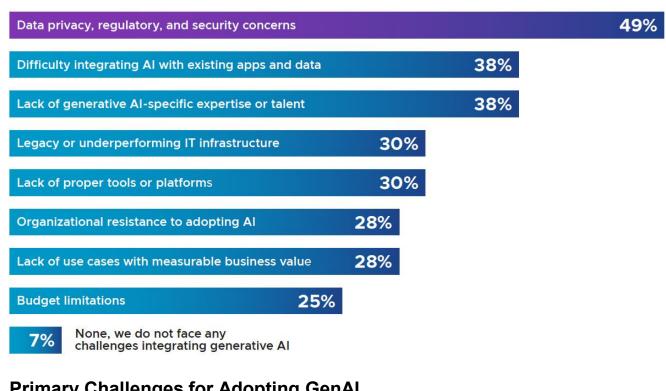
Private Cloud Outlook 2025 - The Cloud Reset

https://www.vmware.com/docs/private-cloud-outlook-2025

The Rise of Generative Al

Organizations are eager to harness GenAl. Only 2% report having no plans to adopt GenAl. The remaining 98% are somewhere on the adoption curve, with 77% already running pilots or live deployments.

Use cases being explored fall evenly across digital assistants, predictive analytics, and customer service initiatives. Still, 28% cite a shortage of use cases with clear, measurable business value as a barrier to GenAl adoption.



Primary Challenges for Adopting GenAl

Environments for Running GenAl 56% 55% 17% Public Cloud Private Cloud **Bare Metal**





PREDICTIVE AI







SPECIALIZED AI MODEL

Impact of Generative AI in the Enterprise

PREDICTIVE AI GENERATIVE AI Marketing Supply Chain Sales **Data Scientists** Customer Research and Legal Operations Development Manufacturing Procurement **Predictive Analytics NATURAL LANGUAGE INTERACTION APIs** \vee \vee \vee **SPECIALIZED AI MODEL**

Human

Resources

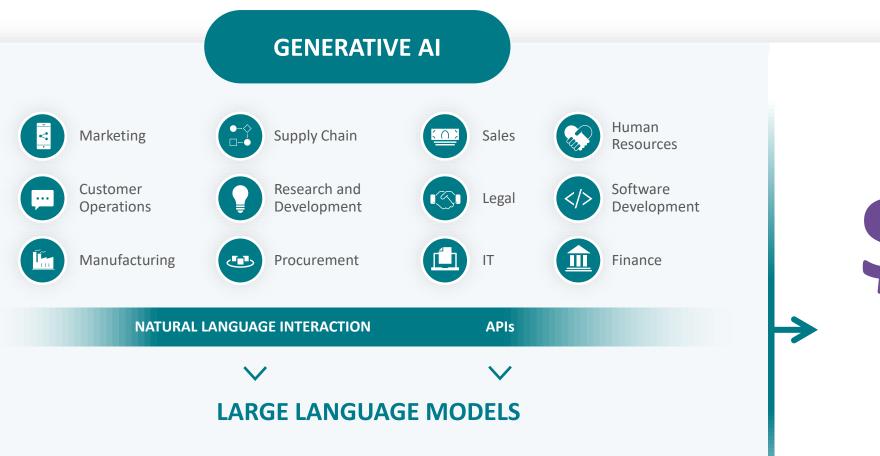
Software

Finance

LARGE LANGUAGE MODELS

Development

Impact of Generative AI in the Enterprise





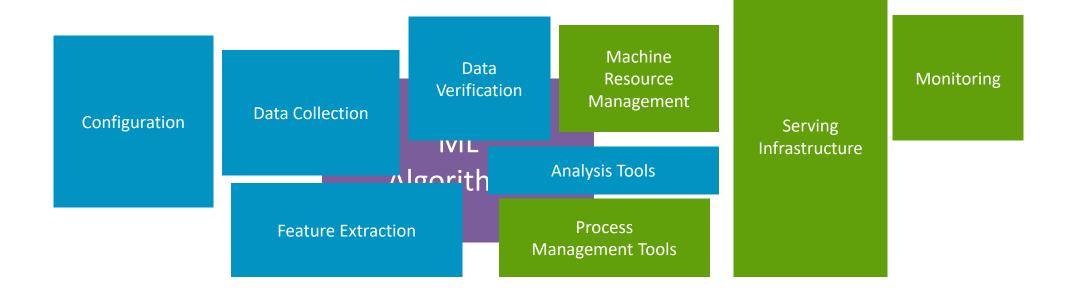
annual economic value

Source: McKinsey, The economic potential of generative AI: The next productivity frontier, June 2023



Complexity of AI Environments

Infrastructure is <u>Half</u> of the Challenge



Corp IT / Infrastructure Domain

Data Scientists / ML Engineer Domain

Based on: Hidden Technical Debt in Machine Learning Systems, Scully, D., et al. https://proceedings.neurips.cc/paper_files/paper/2015/file/86df7dcfd896fcaf2674f757a2463e ba-Paper.pdf





Generative AI Terminology

Key definitions

Traditional or Predictive Al

Generative Al

Large Language Models (LLMs)

Training

Fine-Tuning

Customizing

Retrieval Augmented Generation (RAG)

Inferencing









LLM - Generative AI that is focused on generating text. Examples-GPT-5 (ChatGPT), LLaMA, Claude, Mistral

Training - Educating an AI model from scratch with a lot of data. Requires high costs, massive datasets and high infrastructural resources

Example- Train a LLM from scratch on millions or trillions of data points

Fine-Tuninig - Take a pretrained model and re-train it with your domain data

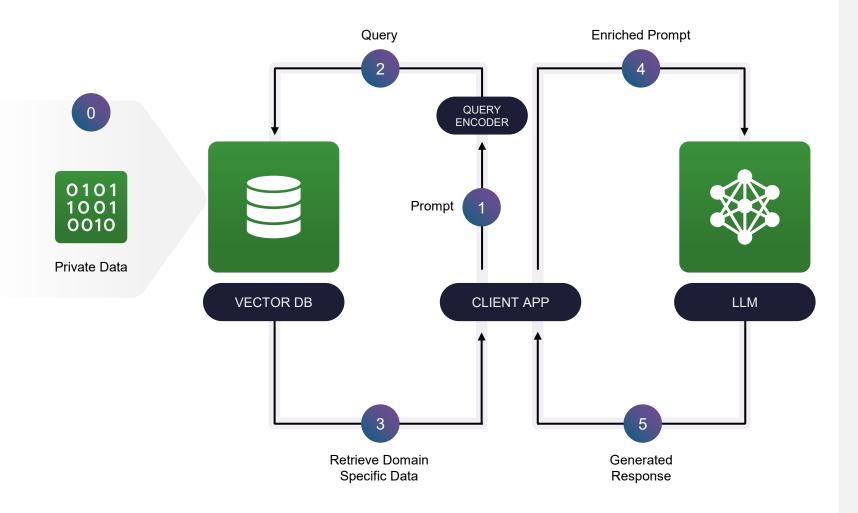
Example- Bank fine tunes an LLM with suspicious activity reports, regulatory filings and internal policy documents to better understand fraud patterns

Customizing - Adapting an Al model for your needs with better instructions, databases and more, without retraining it. e.g. RAG

Inferencing - Deploy a trained AI model to make predictions or generate outputs based on new input data

Retrieval Augmented Generation (RAG)

Intelligence + data



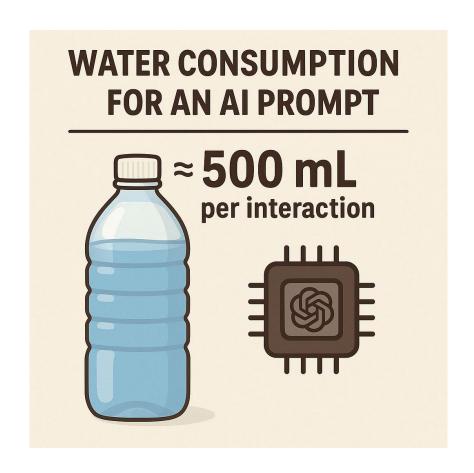
- Data Loading
 Feed your company's private data into the database
- Prompt
 The user generates the prompt
 "How does vSAN ESA handle snapshots?"
- Query
 The prompt is forwarded to the Vector DB which looks for documents, articles, any type of information related to "How does vSAN ESA handle snapshots"
- Retrieve Domain Specific Data
 Information related to "How does vSAN ESA handle snapshots" is sent back to the client application
- The retrieved information is combined with the user prompt and sent as a new, enriched prompt to the Large Language Model (LLM)

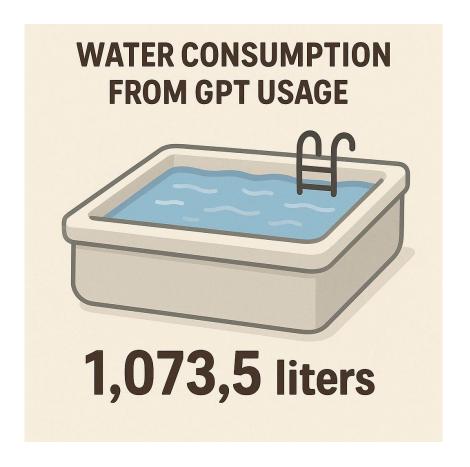
Merging Retrieval with Generation

Generated Response

The LLM uses the enriched prompt (original prompt + retrieved information) and its creative writing skills to create a unique response

We often think of AI as "in the cloud" – but it's deeply physical. It's powered by energy, cooled by water, and grounded in real-world infrastructure





Professor Shaolei Ren at the University of California, Riverside

https://www.researchgate.net/publication/370202417 The Environmental Impact of Al A Case Study of Water Consumption by Chat GPT







VMware Private Al

An architectural approach that balances the business gains from AI with the privacy and compliance needs of the organization.



VCF is now an Al-native platform

Private AI Services are now a part of VCF Subscription

VMware Private AI Foundation with NVIDIA **NVIDIA Foundation Models NVIDIA Optimized Models** Third party & Community Models **OVIDIA** NVIDIA **NVIDIA** NVIDIA **NVIDIA** NeMo™ Retriever NIMTM Microservices Blueprints NIM Operator Al Enterprise VCF Private Al Al Blueprints Agent Builder Model Store Model Runtime Services Quick Start Service Distributed GPU Workload Multitenancy Resource Scheduler Monitoring **VMware Cloud** Security Foundation[™] Hitachi Vantara **D¢LL**Technologies Lenovo Fsas Technologies **Hewlett Packard** SUPERMICE Choice of LLMs **Bare-Metal Performance** Faster Time-To-Value

Why VCF for Private AI? 0



A Day in the Life of a Generative AI App Developer











Model Selection

Data Science Team determines which LLM to utilize specific to their intended AI Use Cases

Setup

Data Scientist requests Al Workstation and loads model of choice into NIM

RAG

Domain specific data uploaded to Vector DB

Iterative prompt process to test usefulness

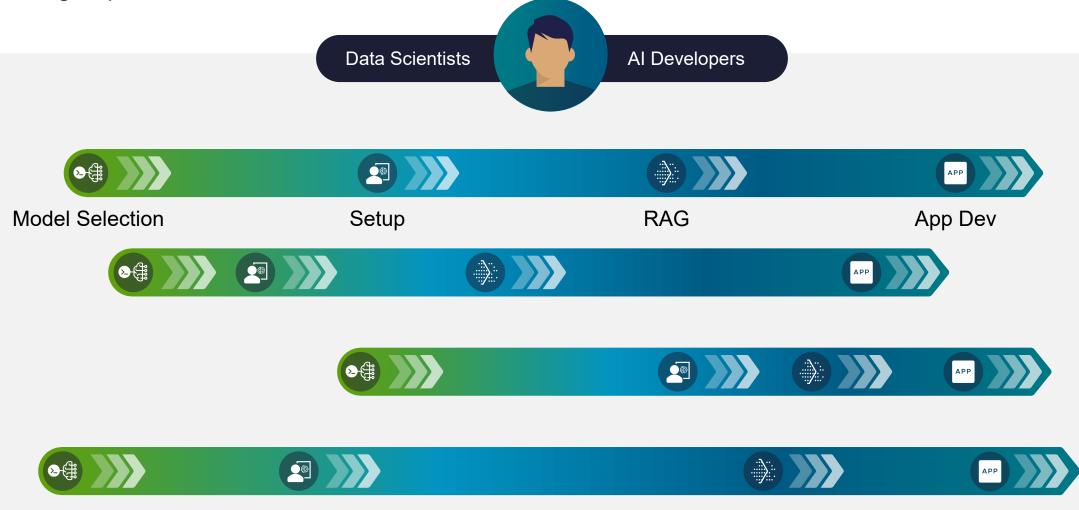
App Dev

NVIDIA NIM API's utilized to develop the end-user's Al application



Teams of Data Scientists

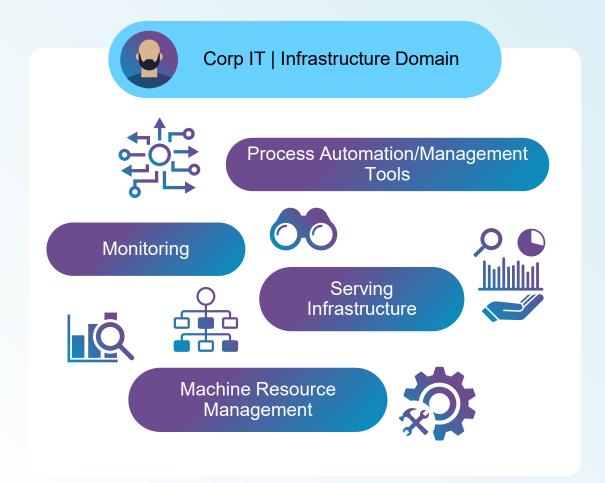
Working in parallel on Models

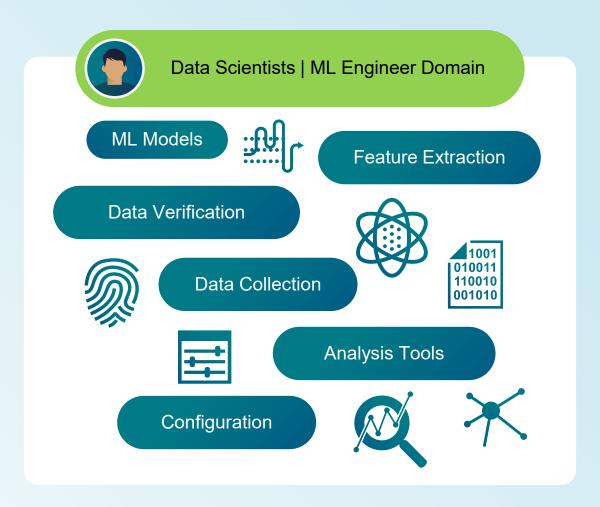




Complexity of AI Environments

Infrastructure is half of the challenge





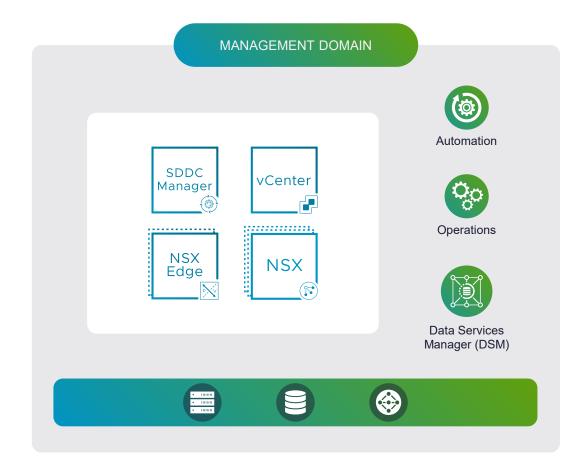
Deployment Options

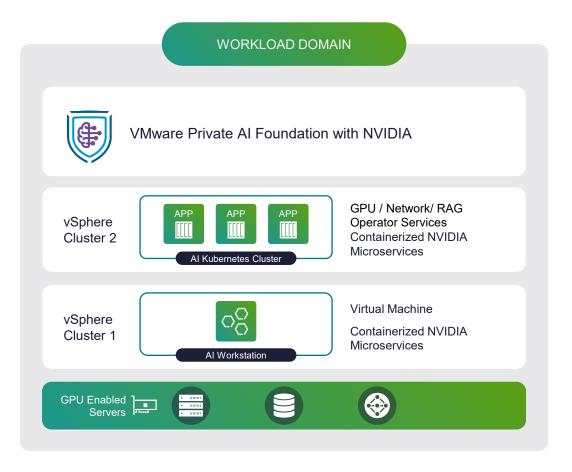




VMware Cloud Foundation

VMware Private AI Foundation with NVIDIA - Add-on VMware components

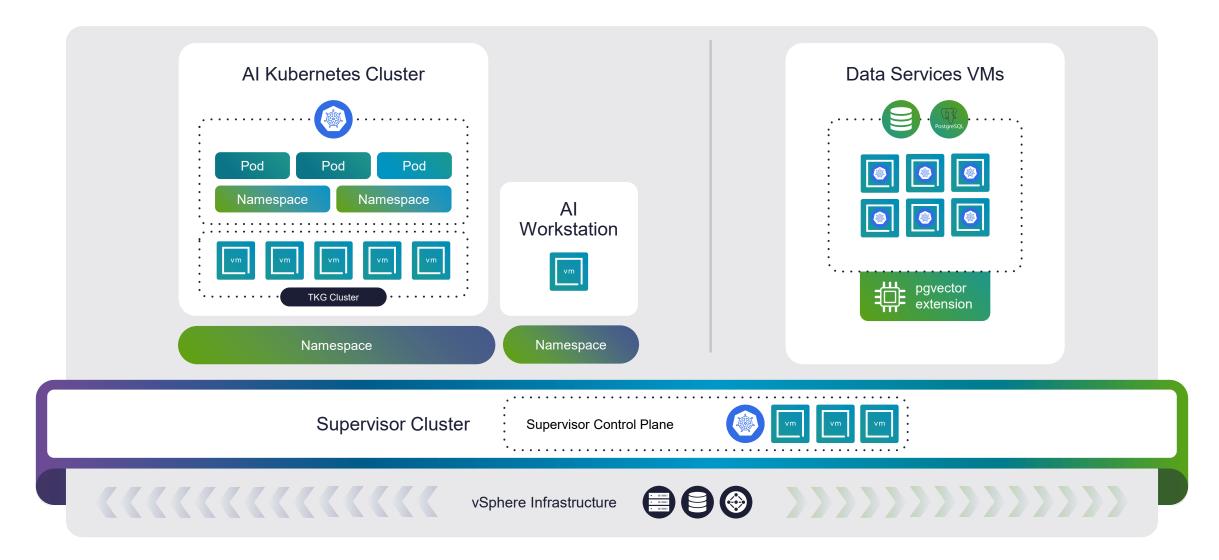




1 laaS Control Plane formerly known as vSphere with Tanzu

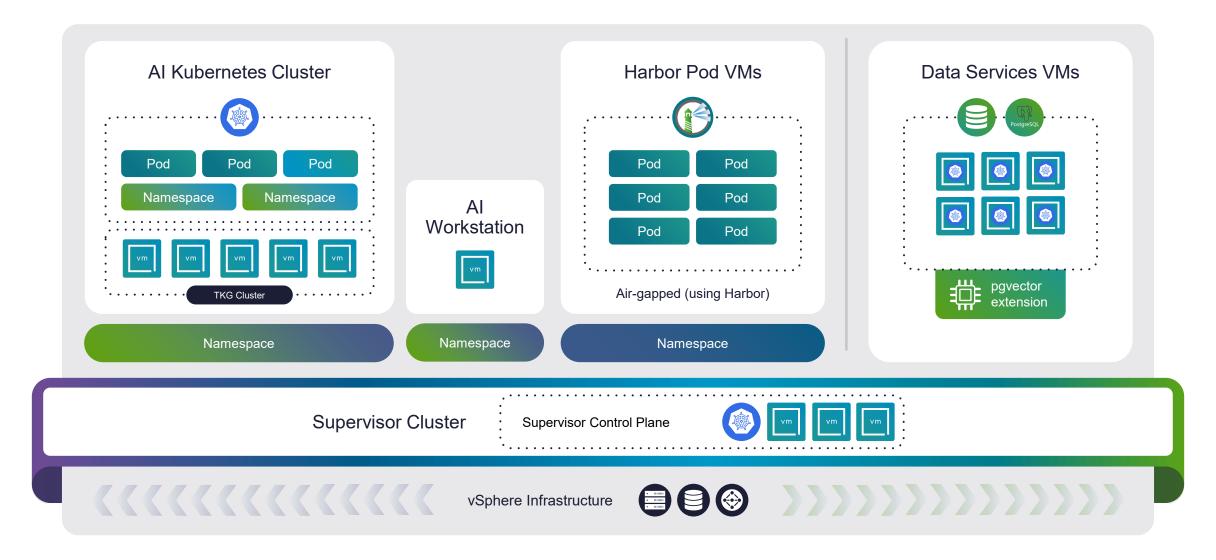


VMware Private AI Foundation with NVIDIA - Deployments





VMware Private AI Foundation with NVIDIA

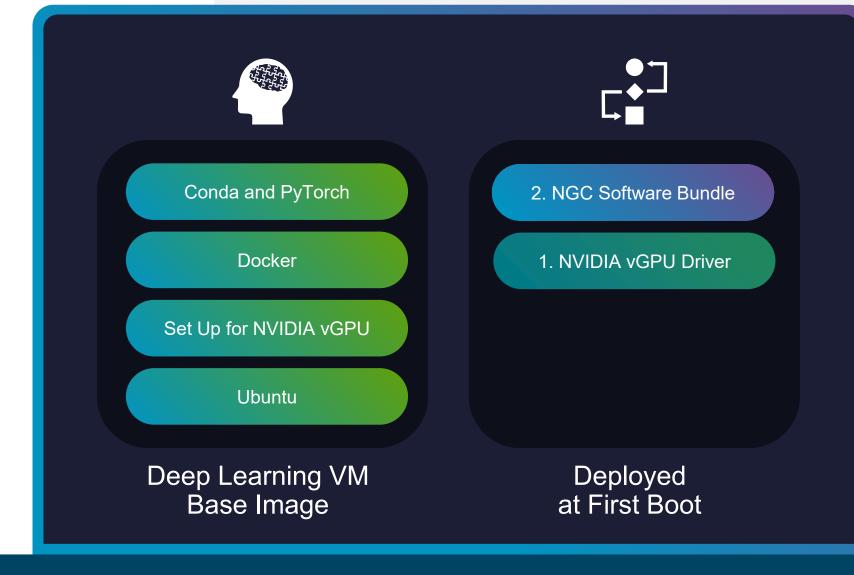




AI Workstation

Deep learning virtual machine



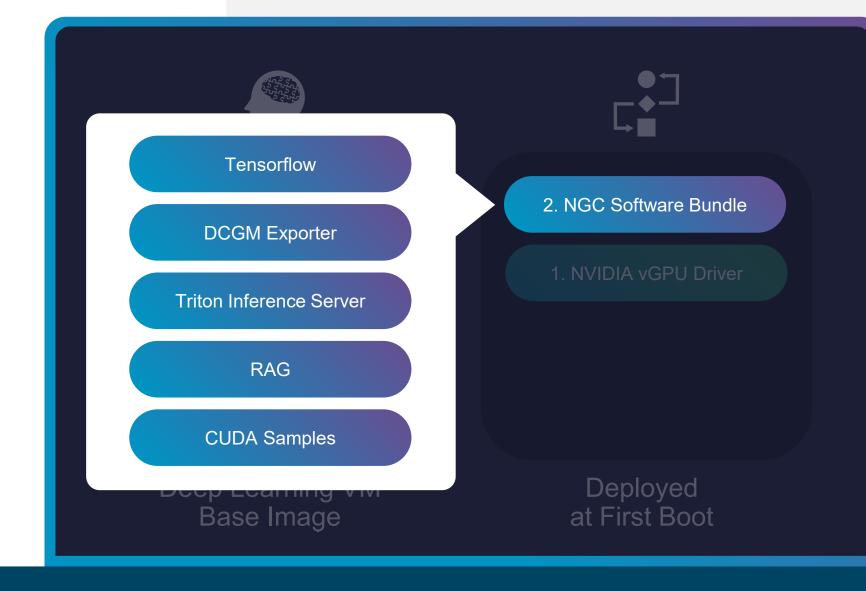




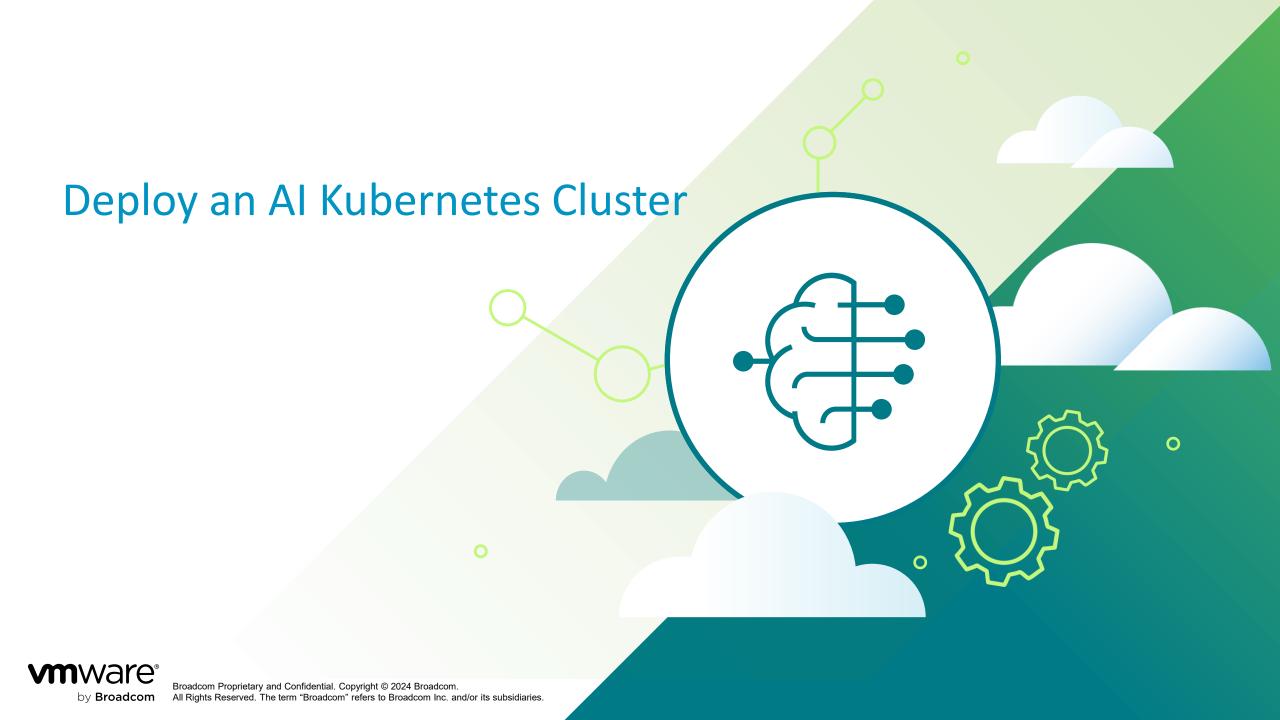
AI Workstation

Deep learning virtual machine

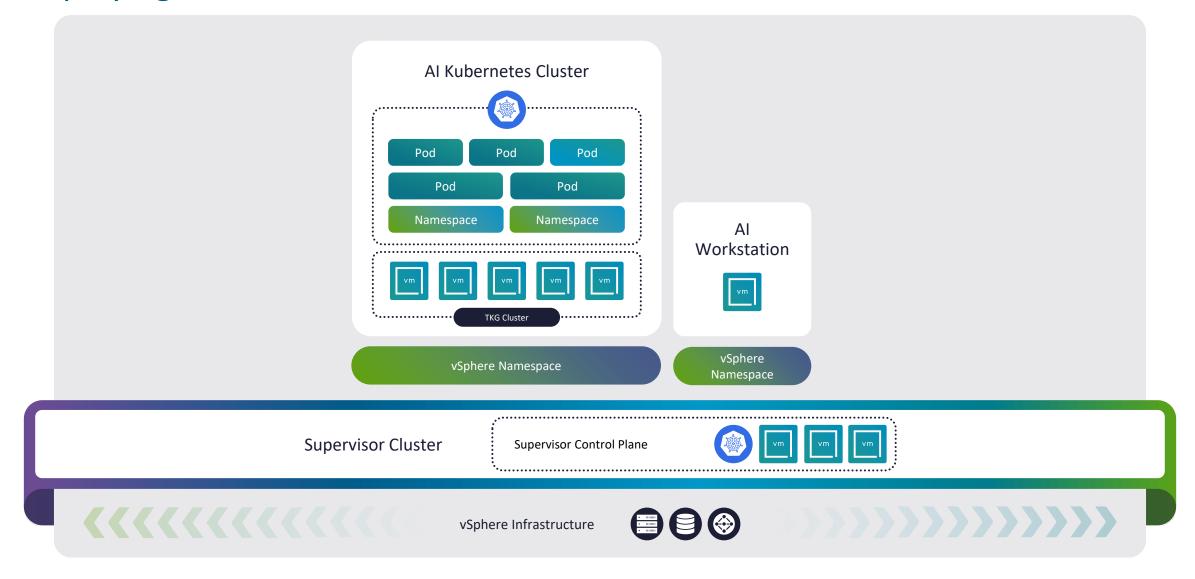






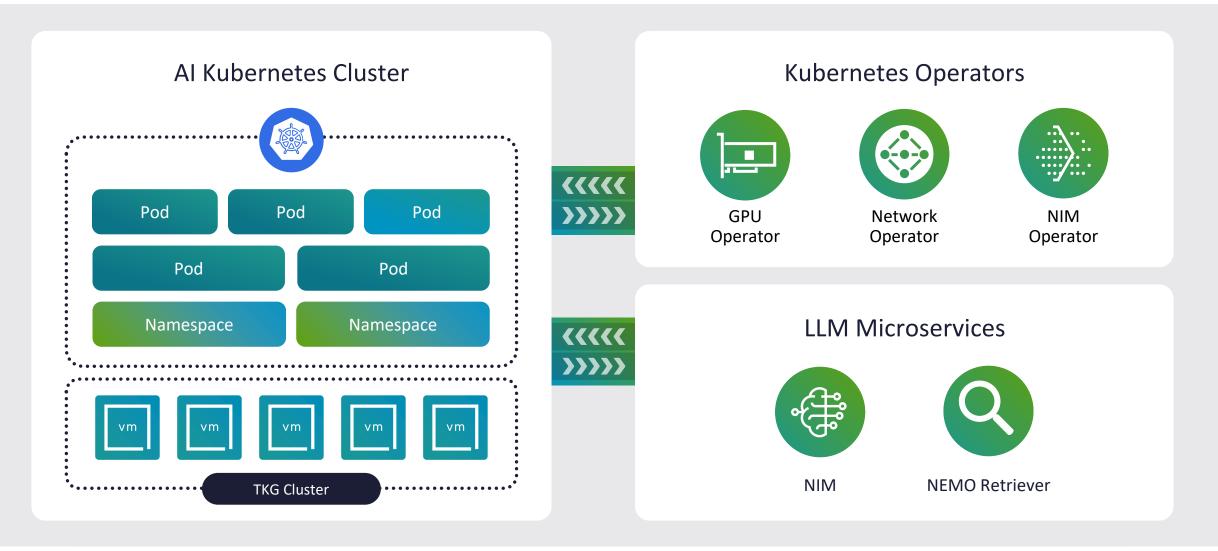


Deploying Microservices



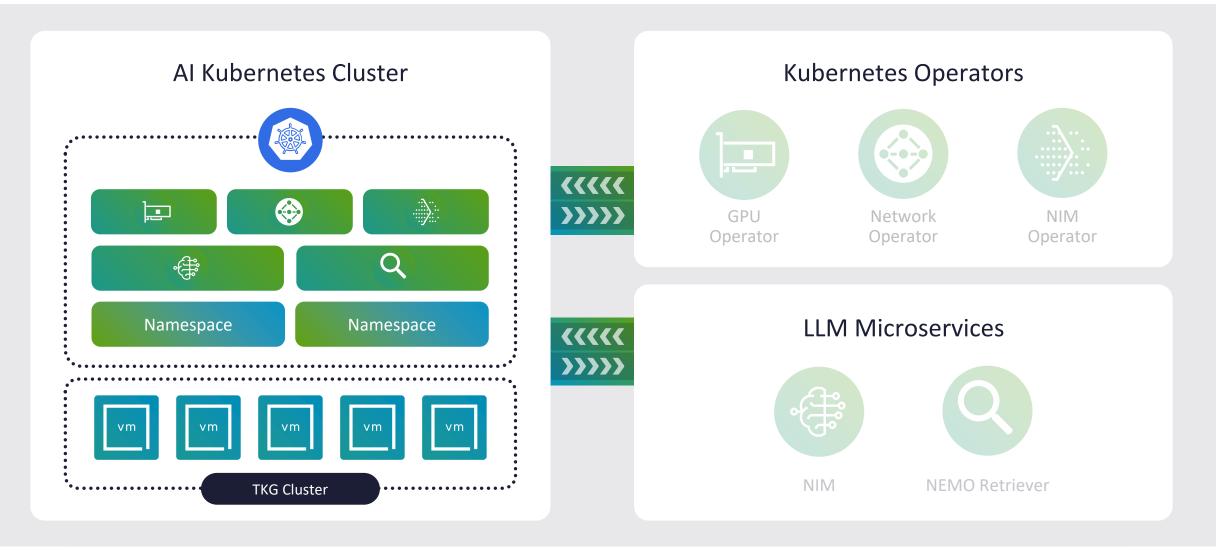


Operators and LLM Microservices (Containers)



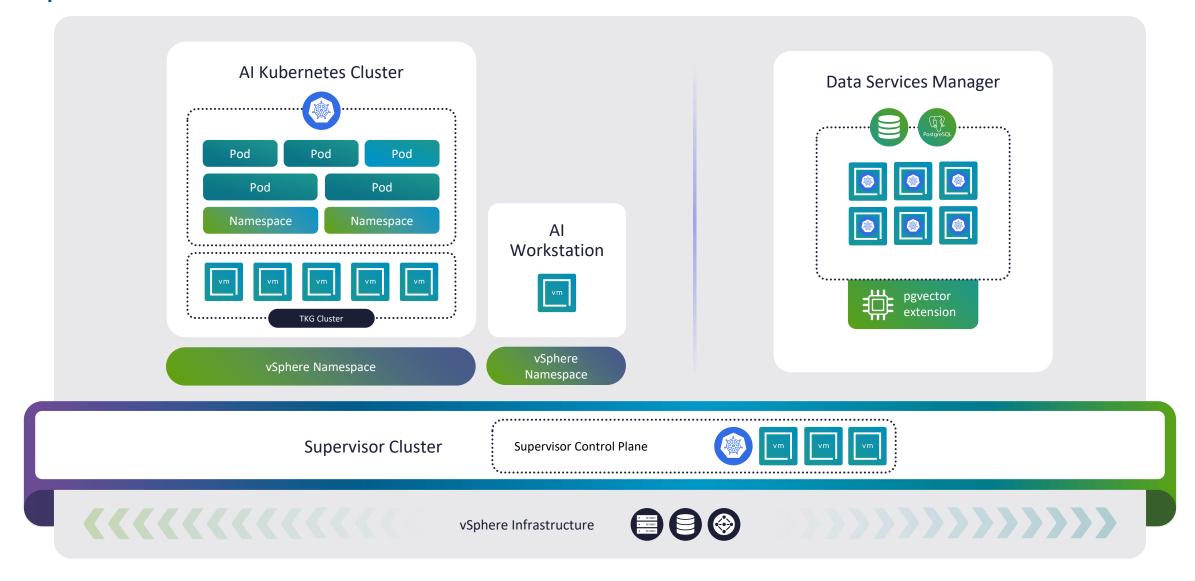


Operator Services



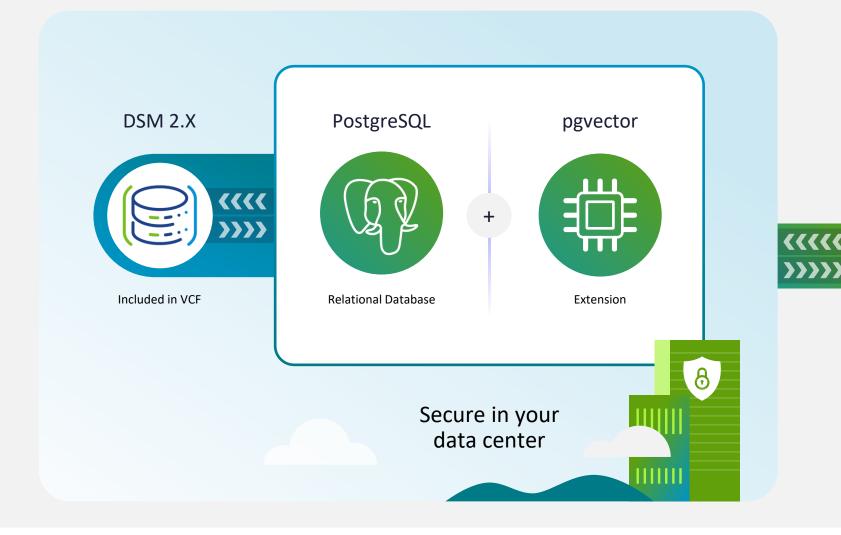


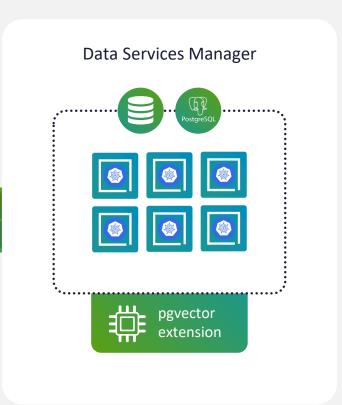
Operator Services





Data Services Manager Deploys PostgreSQL and pgvector







Backup and Availability

Features for vector database backup and replication



Replication scheme with availability



Ensure it's all backed up



Access control with RBAC (role based access control)





GPU Monitoring





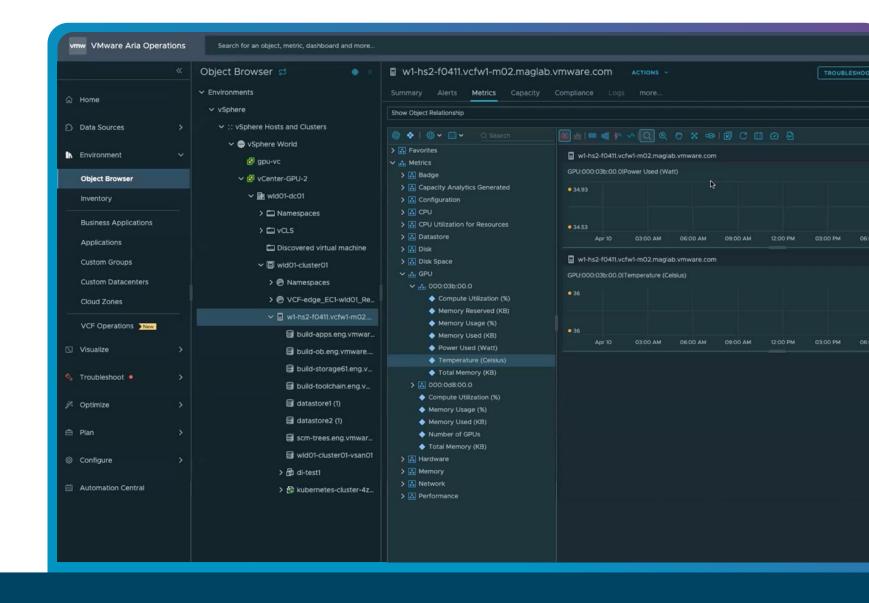
Broadcom Proprietary and Confidential. Copyright © 2024 Broadcom.

All Rights Reserved. The term "Broadcom" refers to Broadcom Inc. and/or its subsidiaries.



GPU Health Monitoring

Enhanced GPU visibility with insight into compute and memory usage across hosts and clusters





Thank you for paraticipating!

Not finished yet ©

November & December
VMware Tech Talks winter edition
Series of afterwork gatherings with hot&cold drinks and hot topics





