# Accelerating Enterprise: Tools and Techniques for Next-Generation AI Deployment

Mahan Salehi, Product Manager, NVIDIA
Nave Algarici, Product Manager, NVIDIA

GTC, March 2024

# Enterprise are on the Generative AI Journey

## 2022

### Explosion

ChatGPT gets announced late in 2022, gaining over 100 million users in just two months. Users of all levels can experience AI and feel the benefits firsthand.

## 2023

### Experimentation

Enterprise application developers kick off POCs for generative AI applications with API services and open models including Llama 2, Mistral, NVIDIA, and others.
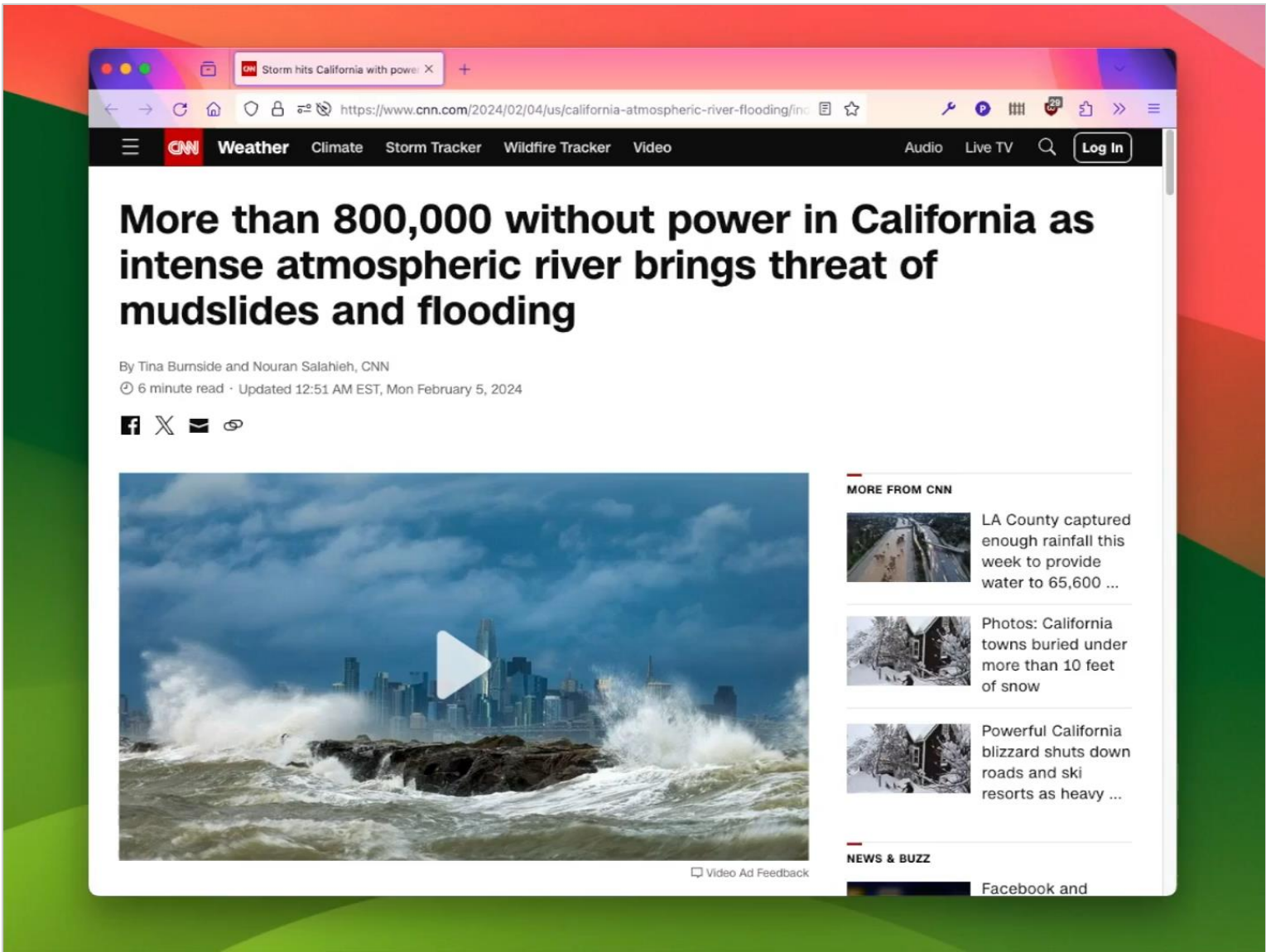
## 2024

### Production
### (aka "Inference")

Organizations have set aside budget and are ramping up efforts to build accelerated infrastructure to support generative AI in production.
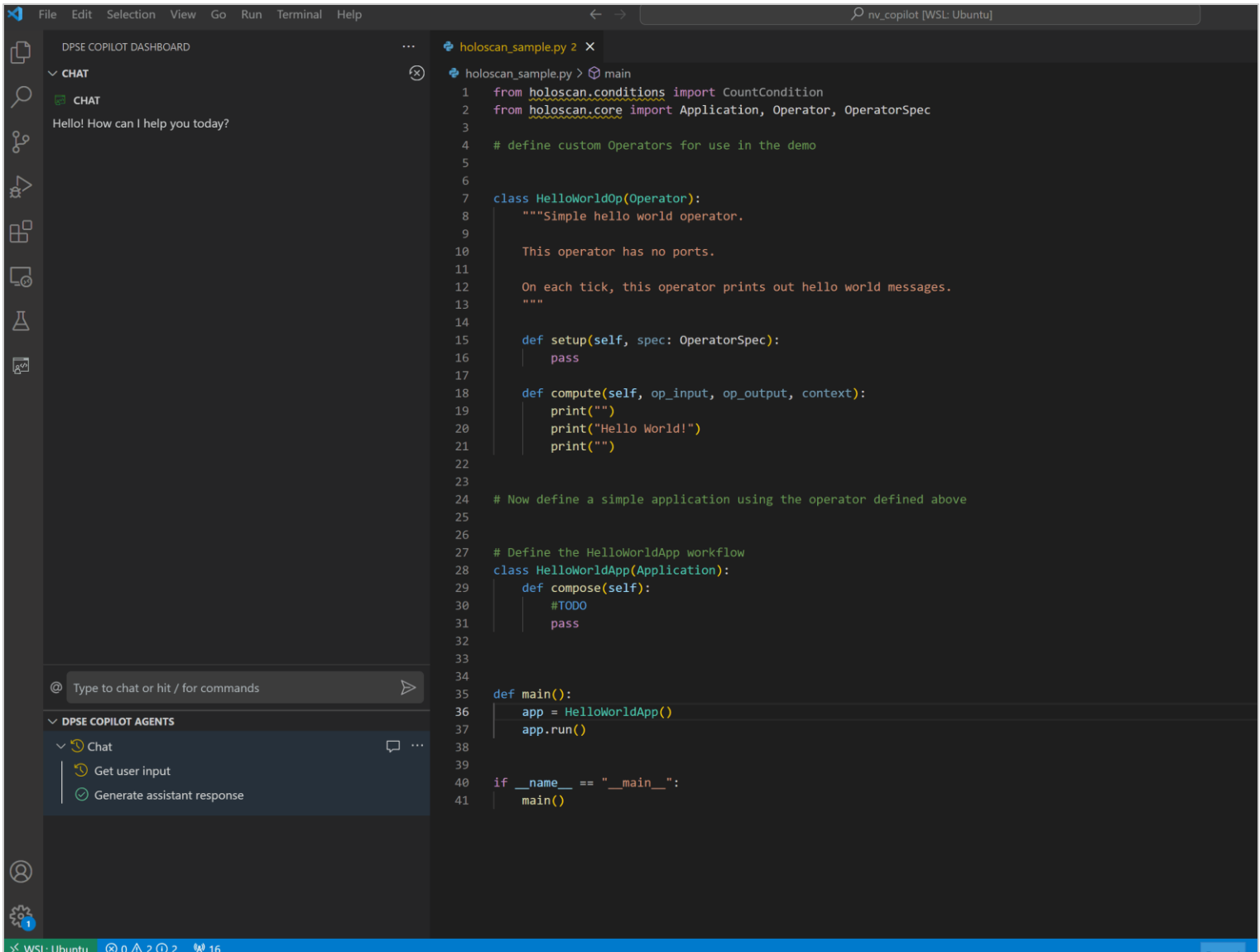
NVIDIA.

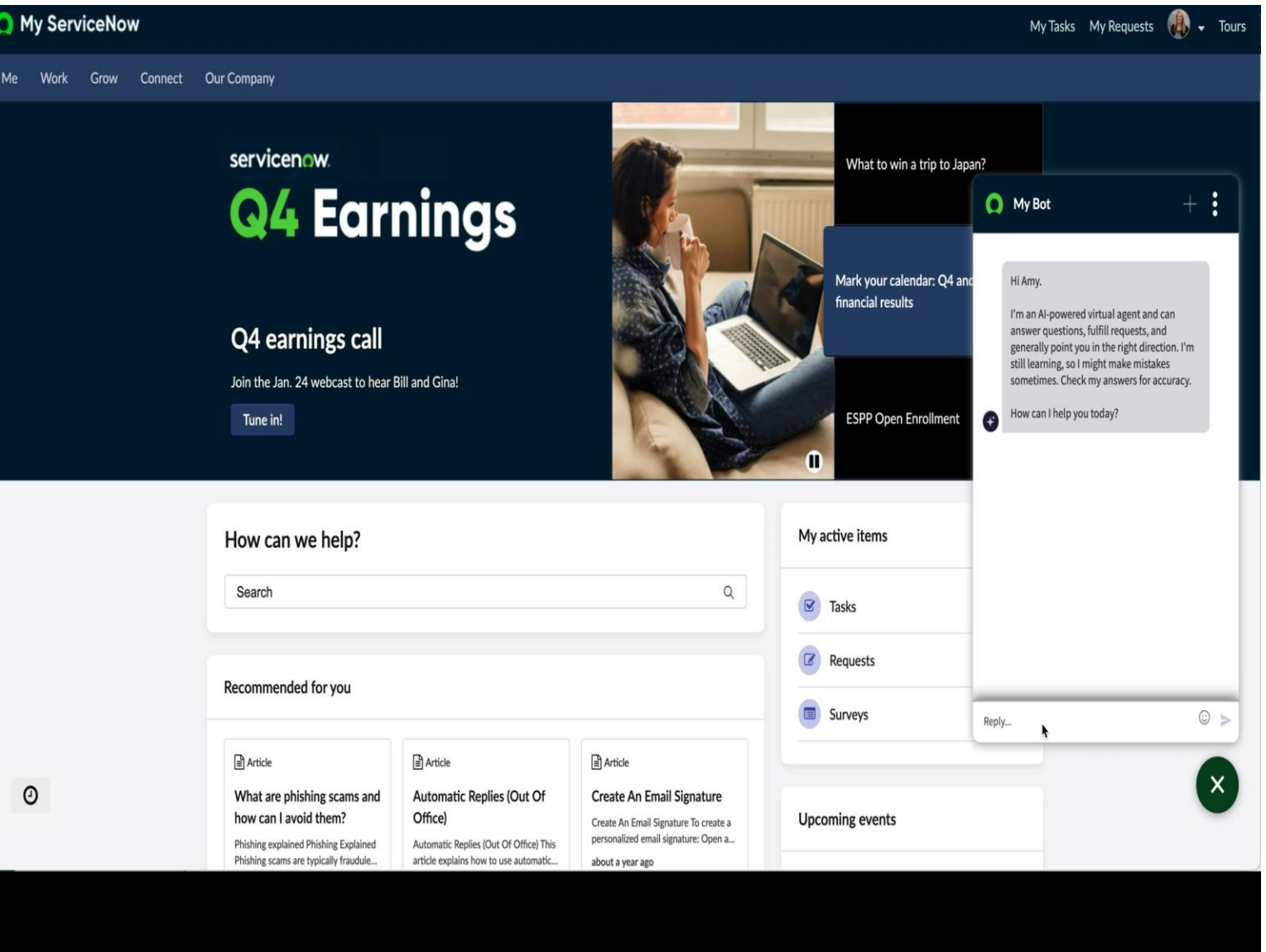# Activate the Potential Within Your Organization

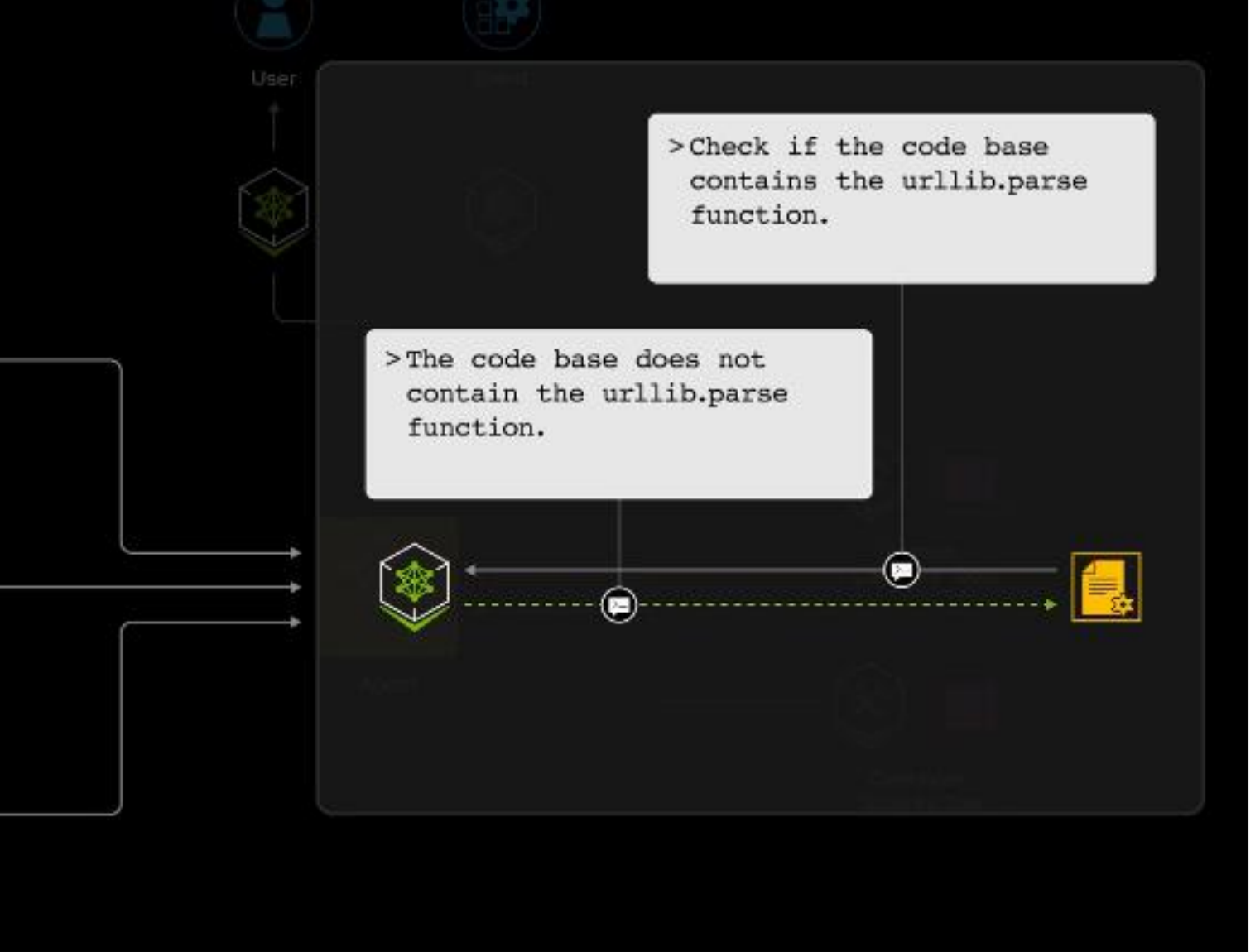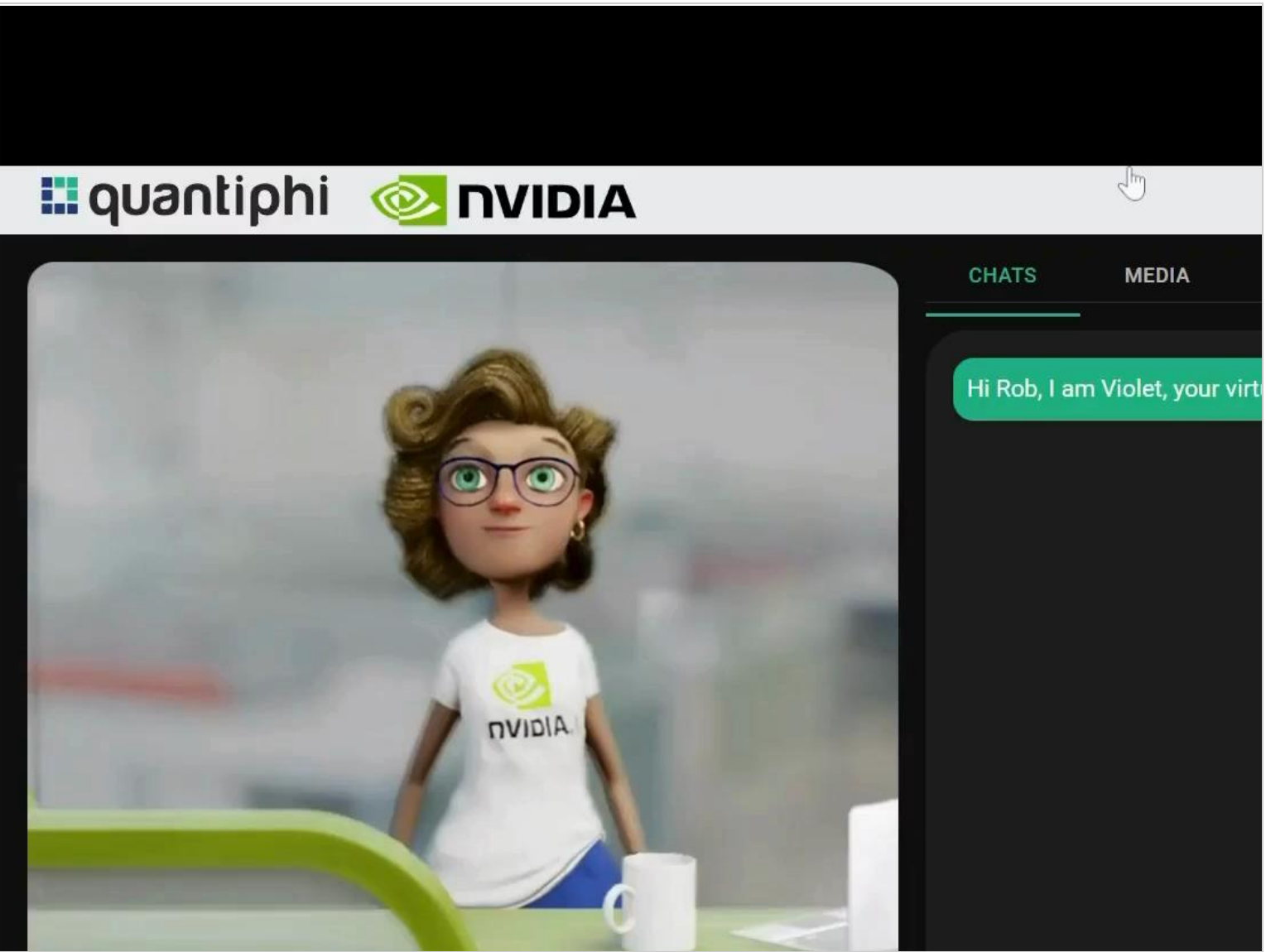Across Every Industry and Job Function



Kinetica Telco Copilot
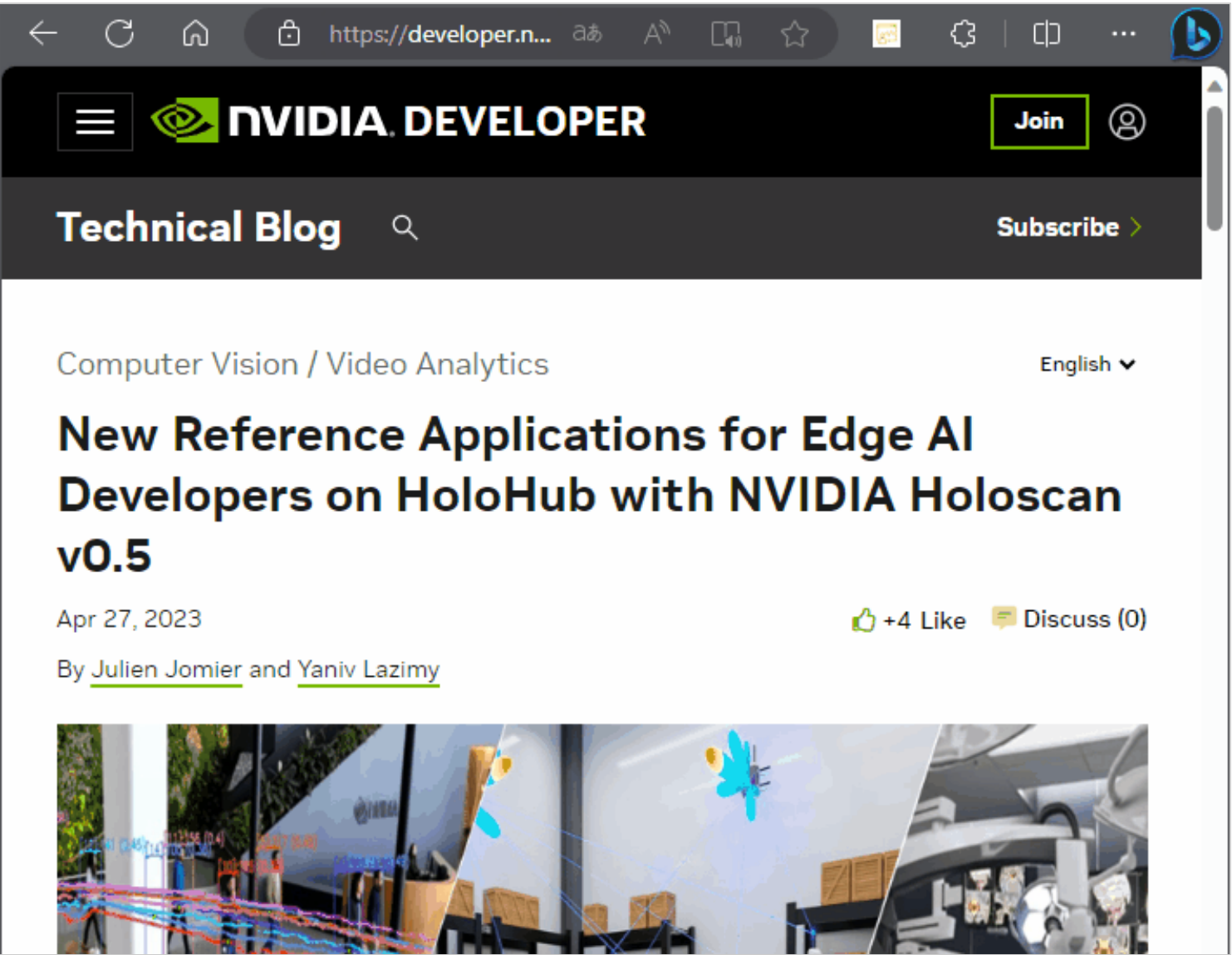


Coding Copilot



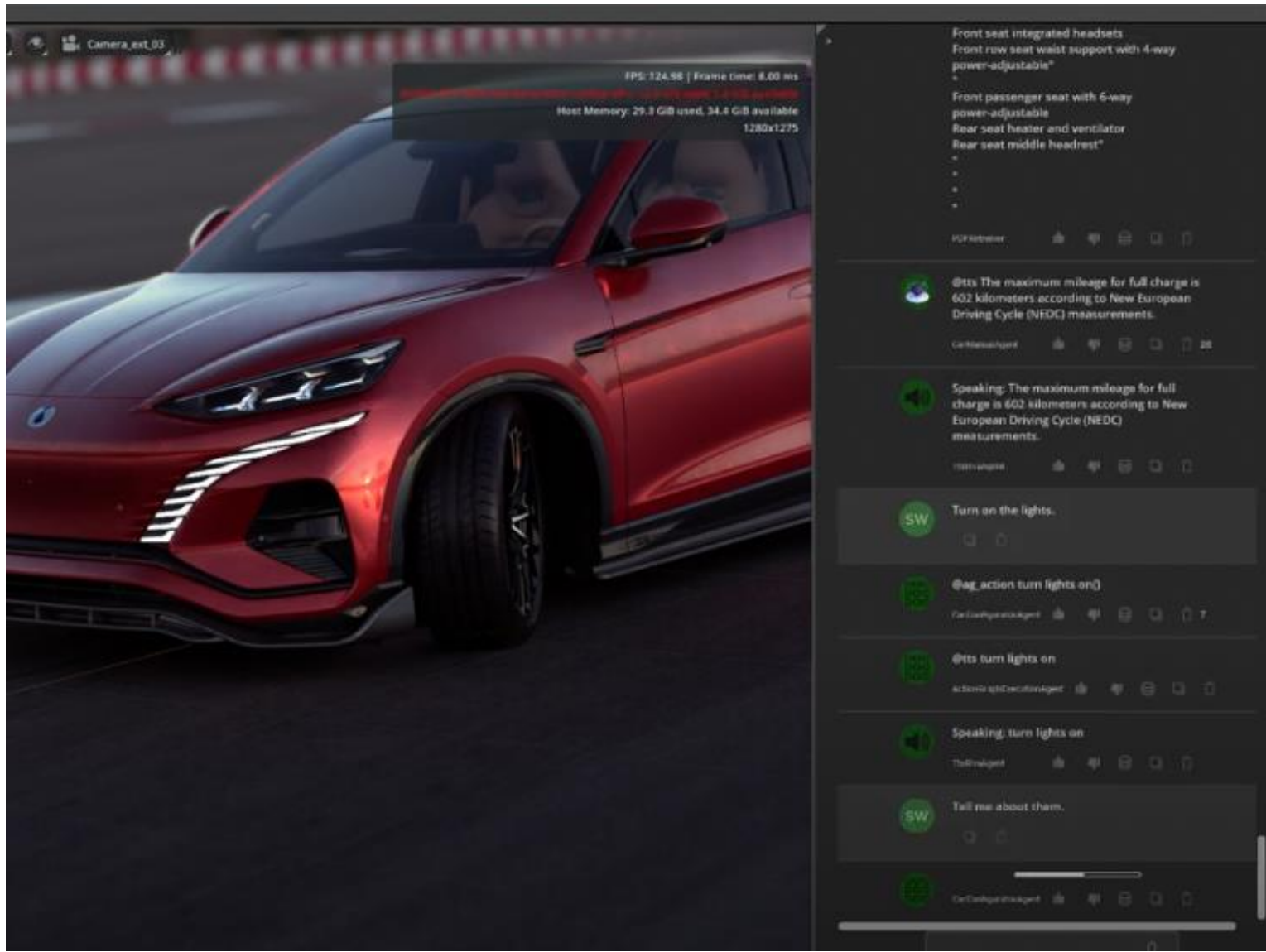ServiceNow Customer Relations Management
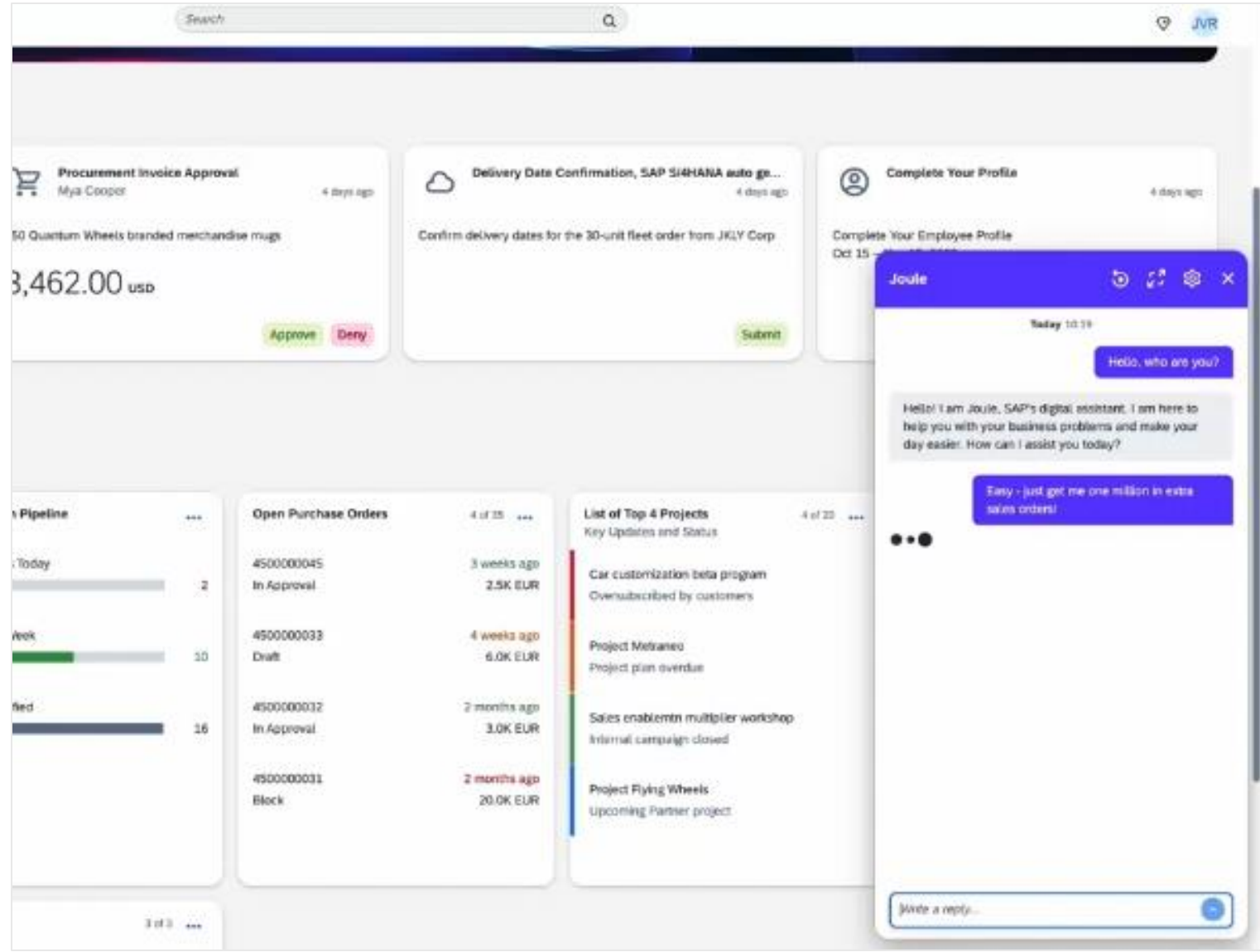


Software Security Analysis



Telco Customer Service Avatar



Summarization



Car Configurator



SAP CRM Assistant

# Generative AI Deployment Options

## Enterprises experimentation with generative AI applications



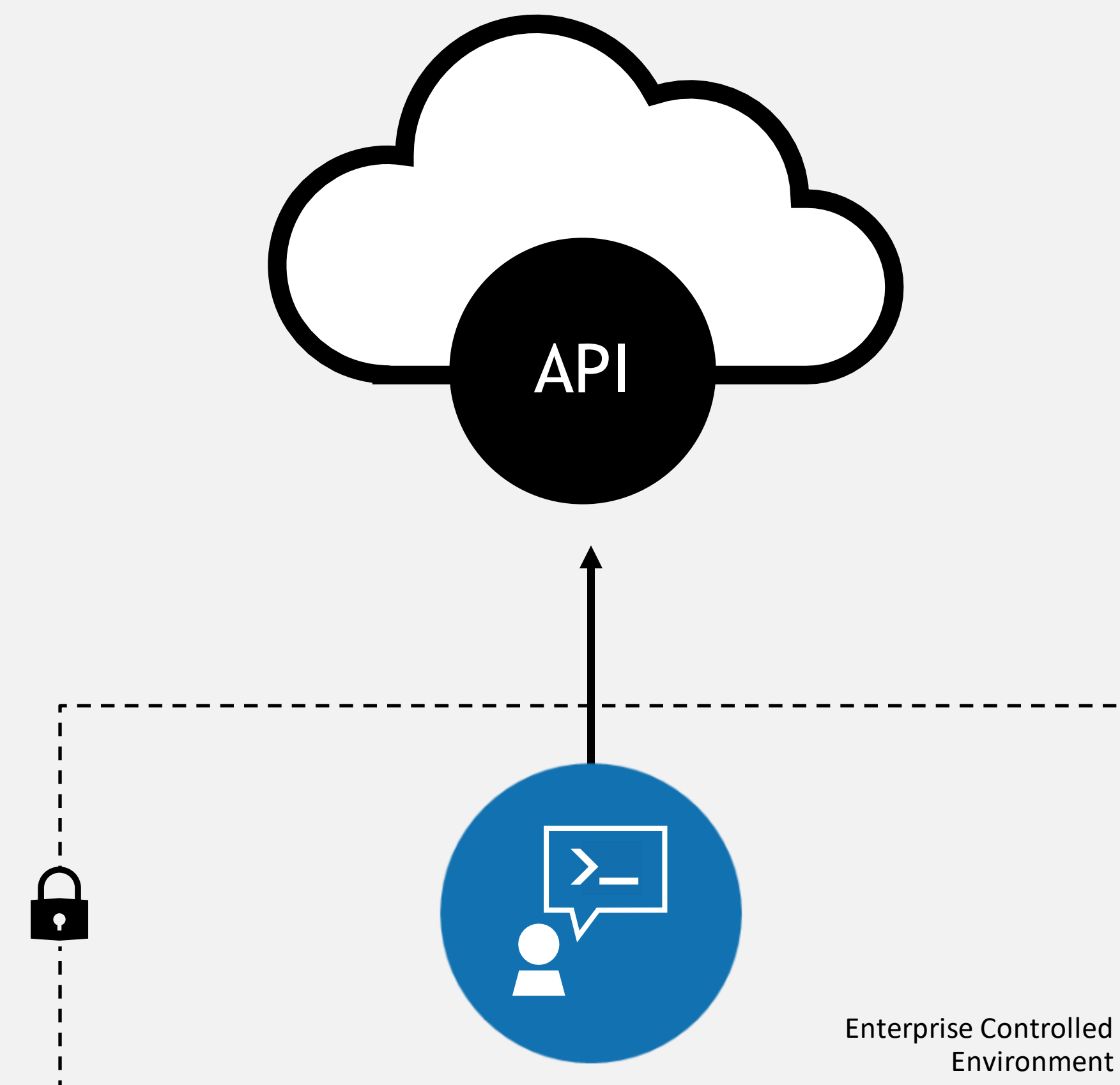**Managed Generative AI Services**

**Open-Source Deployment**

**Easy to use** APIs for development

**Fast path** to getting started with AI

**Infrastructure limited** to managed environment

Data and prompts are **shared externally**

**Limited control** for overall generative AI strategy
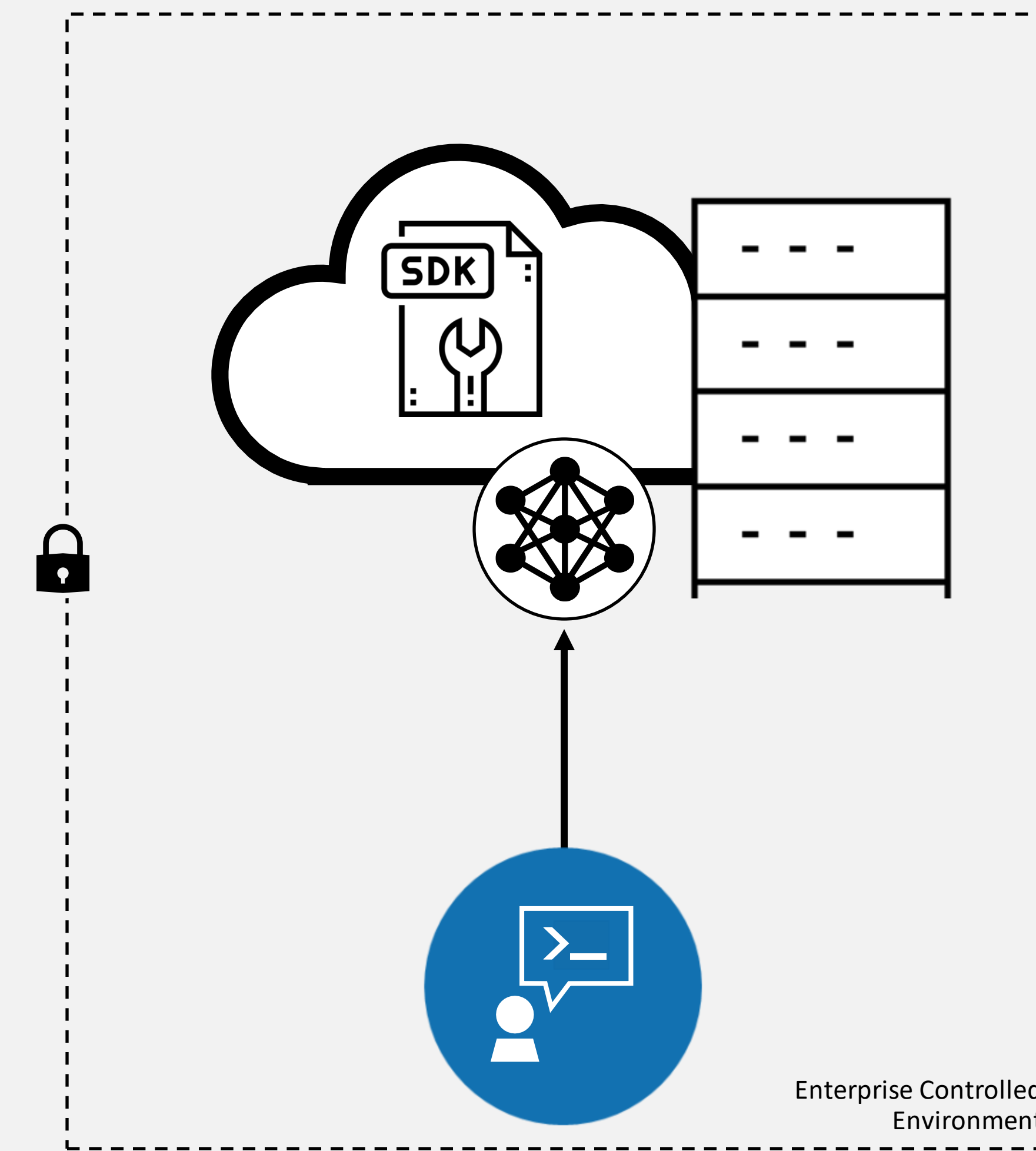
API

Enterprise Controlled Environment

**Run anywhere** across data center and cloud

**Securely manage** data in self hosted environment

**Tuning required** for different infrastructure
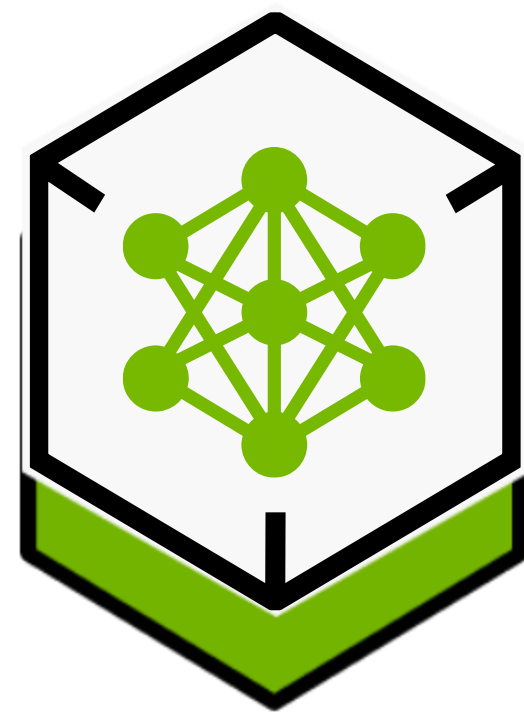
**Custom code** for APIs and fine-tuned models

**Ongoing maintenance** and updates

SDK

Enterprise Controlled Environment
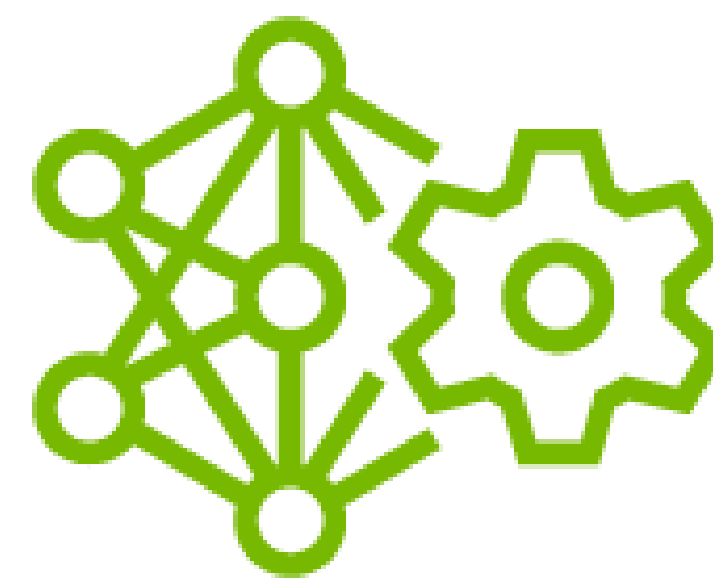
# Inference Microservices for Generative AI

Easiest, fastest, and most portable way to put generative AI models in production. Deploy in 5 minutes.

## Deploy Anywhere

Deploy at scale on preferred infrastructure.
Maintain control of generative AI models and data
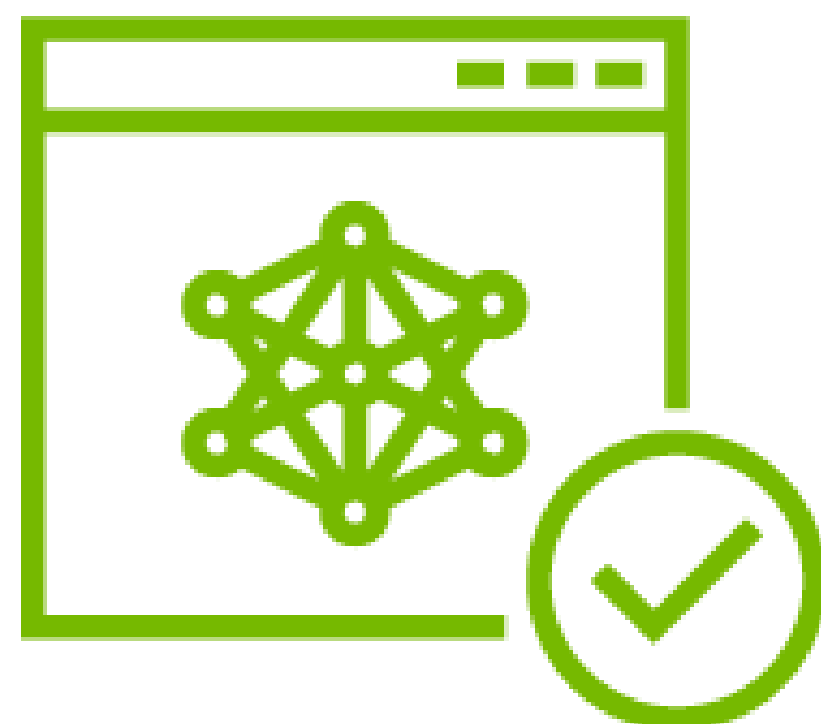
## Easy to use Industry Standard APIs

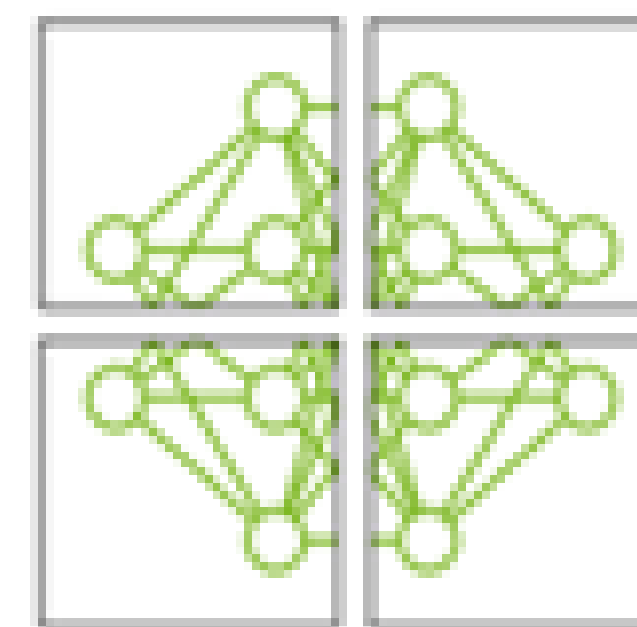Easy to integrate. Supports OpenAI API protocol

## Day 0 Model Support

Domain specific code for each domain, including LLMs, VLMs, video, healthcare, and more

## Improve Cost & Performance
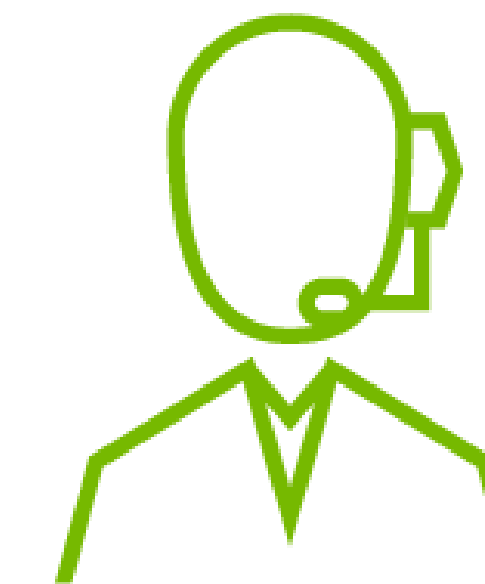
Optimized to provide best latency, throughput, and cost performance

## Supports Tuned Models
## for Best Accuracy

Deploy models that are tuned with proprietary data

## Production Ready

Enterprise support with NVIDIA AI Enterprise

NVIDIA.

# Inference Microservices for Generative AI

Easiest, fastest, and most portable way to put generative AI models in production. Deploy in 5 minutes.

**NVIDIA NIM**

**Prebuilt container and helm chart** tested and validated across infrastructure

**Industry standard APIs**
NVIDIA Cloud standards, OpenAI

**Domain specific code** for each domain e.g. LLMs, VLMs, video, healthcare, and more

**Optimized inference models** for each model architecture and hardware SKU

**Support for customized models** build by users targeted use cases (e.g dynamic LoRA, p tuning)

# NVIDIA NIM is the Fastest Path to AI Inference

## Reduces engineering resources required to deploy optimized, accelerated models

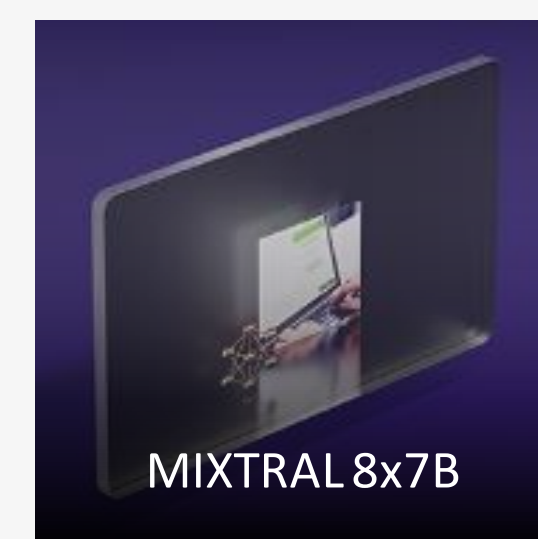| | NVIDIA NIM | Triton + TRT-LLM Opensource |
|---|---|---|
| **Deployment Time** | 5 minutes | ~1 week |
| **API Standardization** | Industry standard protocol OpenAI for LLMs, Google Translate Speech | User creates a shim layer (reducing performance) or   modify Triton to generate custom endpoints |
| **Pre-Built Engine** | Pre-built TRT-LLM engines for NV and community models  MISTRAL AI_   Llama 2   starcoder   NVIDIA Nemotron | User converts checkpoint to TRTLLM format and creates and runs sweeps through different parameters to find the optimal config |
| **Triton Ensemble/  BLS Backend** | Pre-built with TRT-LLM to handle pre/post          processing (tokenization) | User manually sets up + configures |
| **Triton Deployment** | Automated | User manually sets up + configures |
| **Customization** | Supported – P-tuning and LORA, more planned | User needs to create custom logic |
| **Container Validation** | Pre-validated with QA testing | No pre-validation |
| **Support** | NVIDIA AI Enterprise - Security and CVE   scanning/patching and tech support | No enterprise support |

# Inference Microservices for Generative AI

Fastest way to deploy AI models on any accelerated infrastructure across cloud, data center, and PC

## NVIDIA API Catalog



| MIXTRAL 8x7B | GEMMA 7B | FUYU | NEMO RETRIEVER | AI GENERATOR | KOSMOS 2 | 3D GENERATOR | AUDIO2FACE | ESM FOLD | VISTA-3D | DIFFDOCK | MolMIM |

MISTRAL AI_ · Google · ADEPT · NVIDIA · gettyimages · Microsoft · shutterstock · NVIDIA · Meta · NVIDIA · MIT · NVIDIA

## NVIDIA NIM



Microsoft Azure · aws · Google Cloud · ORACLE · DELL Technologies · Hewlett Packard Enterprise · Lenovo · SUPERMICRO

NVIDIA

# Journey to NVIDIA Cloud APIs
## NVIDIA API Catalog - Taking our NIMs to Market

**1** Landing Pages

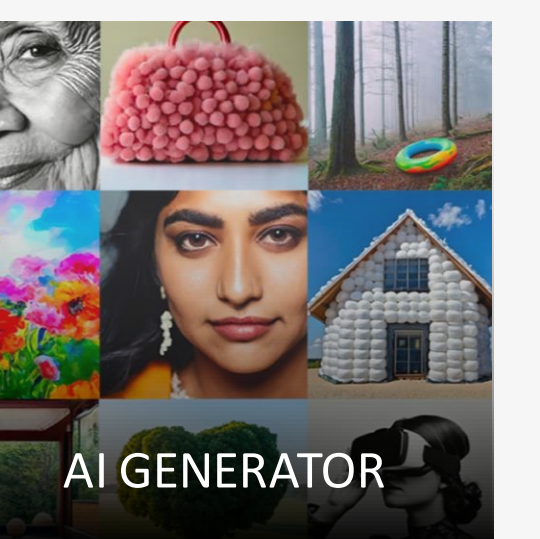*Part of campaign strategy*

Platform Pages



NVIDIA AI

NVIDIA Omniverse

Industry Pages

Healthcare

Robotics

**2** NVIDIA API Catalog

*Includes all APIs across NVIDIA*

**3** Download & self host API

**3** Route to Partner

## Discover APIs

## Try & develop with hosted APIs

## Consume APIs

# The Amount of Enterprise Data is Massive & Growing

NVIDIA Accelerated Retrieval-Augmented Generation

**20 ZB**

Unique Enterprise
Data Created

2027

**11 ZB**

2024

**800,000** Libraries of Congress

**83%** unstructured data

**50%** audio and video

Source: IDC Global DataSphere

NVIDIA

# What is Retrieval-Augmented Generation?

## Enhance the Accuracy & Reliability of Generative AI

# RAG Challenges

## Difficult to Take a RAG Pipeline from PoC to Production

| Accuracy | Data Security | Complexity | Cost | Innovation Velocity |
|---|---|---|---|---|
| Accurate generations require retrieval systems tuned to match the data | Sending sensitive data to remote endpoints is inherently insecure | System builders must piece together and integrate many components | Automated LLM transactions make transaction costs unpredictable | New models and techniques appear every day |

# NVIDIA Retriever for World-Class Information Retrieval

## Lowest Latency, Highest Throughput, Maximum Data Privacy

### Optimized Inference Engines

Built on NVIDIA TensorRT & Triton
Tuned for RAG application workloads

### World Class Models & Community Model support

SOTA commercial models
Converts to TensorRT for optimal performance

### Flexible & Modular Deployment

Deploy at scale on preferred infrastructure

### Customizable Pipelines & Models

Compose pipelines with microservices
Customize models for target domain

### Production Ready

Enterprise support running in any cloud or on prem

NVIDIA.

# NeMo Retriever Supercharges RAG Applications

World Class Accuracy and Throughput



**2X** World-class accuracy with nearly 2x fewer incorrect answers

**7X** Faster embedding inference throughput

Optimized Inference Engines

World class models and community model support

Flexible and modular deployment

Customizable models and pipelines

Production Ready

# NVIDIA NeMo Retriever: Supercharges your RAG Application

## Enterprise-ready microservices for RAG

Supercharge Software Delivery With Event-Driven RAG

# World Class Accuracy and Throughput

### Retrieval, embedding, reranking microservices for RAG



**Nearly 2X Fewer Incorrect Answers**

- Lexical Search (BM-25): 24%
- E5 Unsupervised Embedding: 48%
- NVIDIA E5 Embedding + Mistral Reranker: 72%

**7X Increase in E5-Large Embedding Inference Throughput**

- A100: 1X
- A100 with NeMo Retriever Embedding microservice: 3X
- H100 with NeMo Retriever Embedding microservice: 7X

Comparing NVIDIA Text QA Embedding Model vs Other Available Options. Recall Top 5, 300 token chunk size, averaging across representative customer datasets from Telco, IT, Consulting, Energy

NVIDIA.

# Simplifying Retrieval at Scale
## NeMo Retriever in Practice

## Download models

The expected models can be seen in the `command`'s of the `docker-compose-ea.yaml` services. We require two models to run this example, an embedding model and reranking model. The state of the art NVIDIA Retrieval QA Embedding model and QA Reranking model.

```
ngc registry model download-version --dest models "ohlfw0olaadg/ea-participants/nv-embed-qa:4"
ngc registry model download-version --dest models "ohlfw0olaadg/ea-participants/nv-rerank-qa-mistral-4b:1_A100"
chmod -R o+rX models # updating read permissions to ensure container can read mounted models directory
```

## Start the sevices

```
docker compose -f docker-compose-ea.yaml up
```

**0** Set up the retrieval pipeline in just a few lines

**1** Create document collection

### Create a collection using the Create Collection endpoint

```
curl "http://localhost:1984/v1/collections?pretty=true" \
  -H 'Content-Type: application/json' \
  -d '{"name": "my collection", "pipeline": "hybrid"}'
```

You'll get a collection ID back in the response which you will use later on. Save it with:

```
export COLLECTION_ID=id-from-the-above-resonse
```

**2** Upload documents

### Add documents to the collection

```
curl "http://localhost:1984/v1/collections/$COLLECTION_ID/documents?pretty=true" \
  -H 'Content-Type: application/json' \
  -d '[
    {
      "content": "This is some text that we are going to put into our index",
      "format": "txt",
      "metadata": {
        "filename": "my-small-file.txt"
      }
    }
  ]'
```

Note: TXT file uploads are limited to 5MiB. PDF file uploads are limited to 50MiB. Collections currently have no limits.

### Query your collection

**3** Retrieve relevant data

```
curl "http://localhost:1984/v1/collections/$COLLECTION_ID/search?pretty=true" \
  -H 'Content-Type: application/json' \
  -d '{"query": "please return docs"}'
```

# Case Study: Cadence Design Systems

## 3.3x fewer incorrect answers retrieving from technical documentation

| Recall | Top 1 | Top 3 | Top 5 | Top 10 |
|---|---|---|---|---|
| **Reference Pipeline** | 36% | 52% | 57% | 64% |
| **NeMo Retriever Hybrid Search** | 57% | 70% | 77% | 80% |
| **NeMo Retriever Hybrid Search + Reranker** | 69% | 81% | 86% | 89% |
| **Improvement Factor** | **2x** | **2.5x** | **3x** | **3.3x** |

cadence®

NVIDIA.

# Unlock Petabytes of Enterprise Data
## Transform Data into Business Insights

**Adobe**

Adobe's proprietary AI will help unlock the knowledge inside the world's more than 3 trillion PDFs worldwide.

**CLOUDERA**

Cloudera will expand its generative AI capabilities by integrating NeMo Retriever with Cloudera Machine Learning to unlock the potential of 25 exabytes of enterprise data.

**COHESITY**

Cohesity data platform customers can add generative AI intelligence to their data backups and archives.

**DATASTAX**

Datastax leverages NeMo Retriever and NVIDIA NIM improving performance of RAG applications. Using NVIDIA H100 GPUs they achieve an embedding and indexing latency of 10 ms.

**NetApp**

NetApp unlocks exabytes of data empowering customers to securely "talk to their data" to access business insights.

**PURESTORAGE®**

Pure accelerates time to insight for enterprises using their own internal data for AI training, ensuring the use of their latest data and eliminating the need for constant retraining of LLMs.

**SAP**

SAP plans to add RAG capabilities that enable generative AI applications to more securely access data running on SAP software to improve accuracy and insights, using NeMo Retriever.

**snowflake**

Snowflake customers will be able to utilize NeMo Retriever directly on their proprietary data in the Data Cloud, all while maintaining data security, privacy, and governance seamlessly through Snowflake's built-in capabilities.

# Next Generation of Enterprise Applications Connect LLMs to Enterprise Data

## Retrieval Augmented Generation Improves LLM Performance and Efficiency

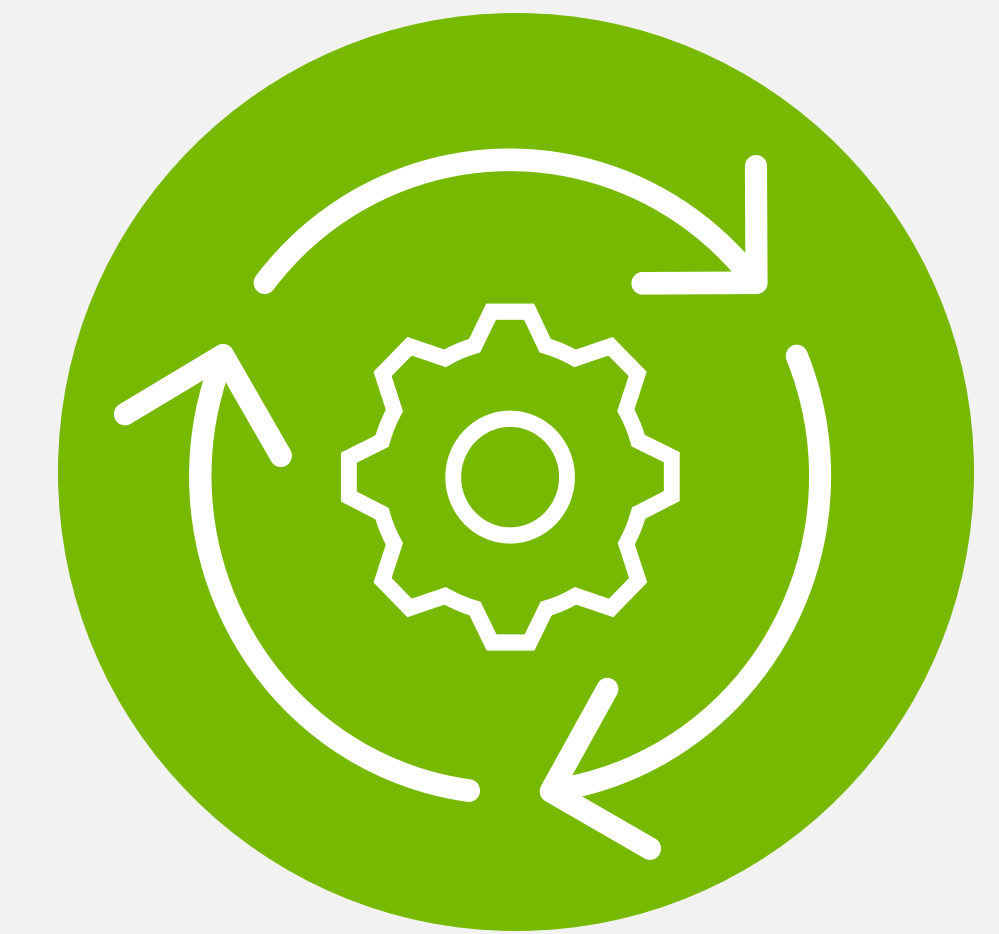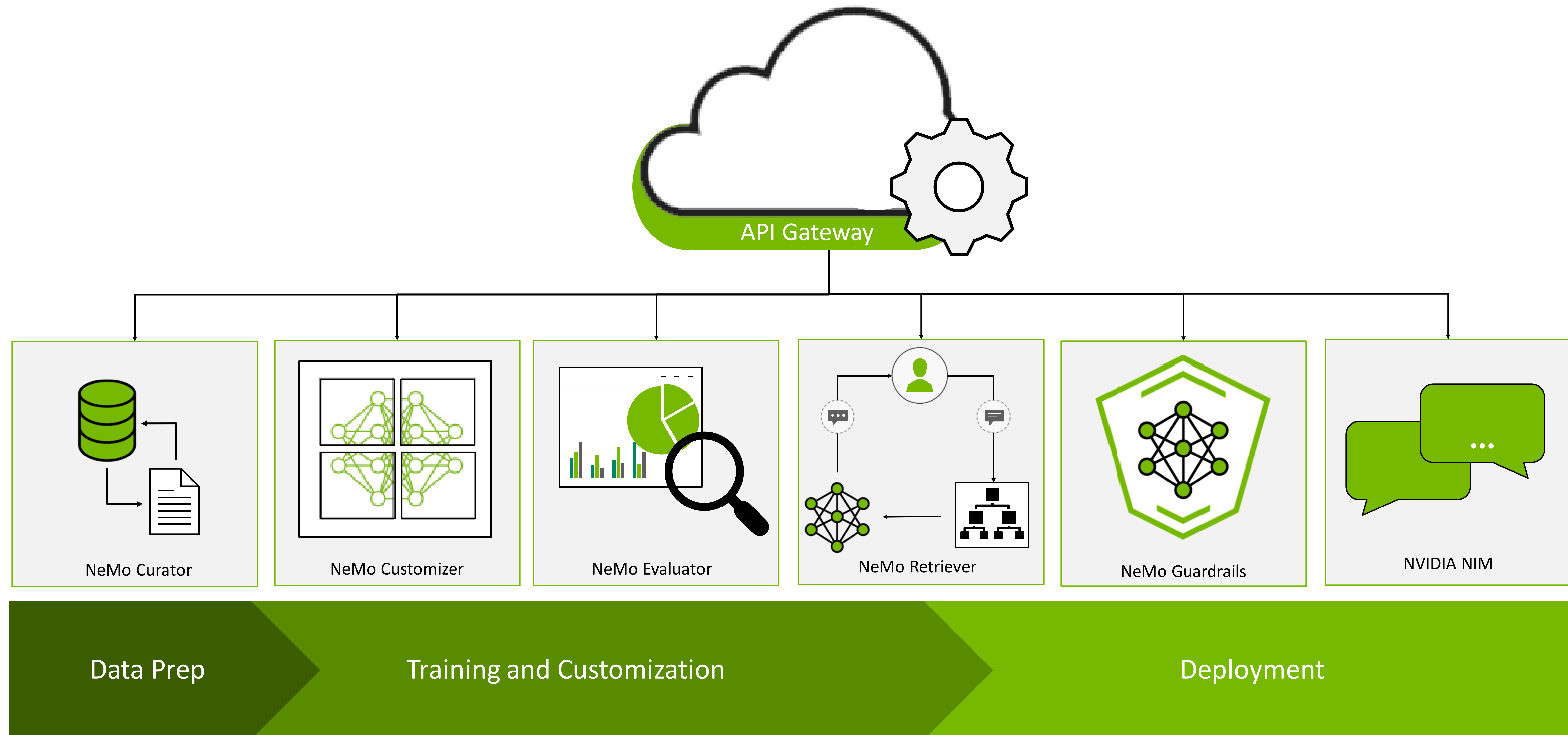| Improved Accuracy | Natural Language Interface | Contextual Understanding | Reduced Computational Costs | Improved Efficiency |
|---|---|---|---|---|
| Models can answer questions about information without having been trained on that data | Human-readable output texts that are easier for people to understand, raising user trust | AI models better understand context when generating text or other outputs | Reduced computational costs from retraining and model size at inference | Models can produce diverse outputs without sacrificing accuracy or efficiency |

NVIDIA.

Building Generative AI Applications for the Enterprise

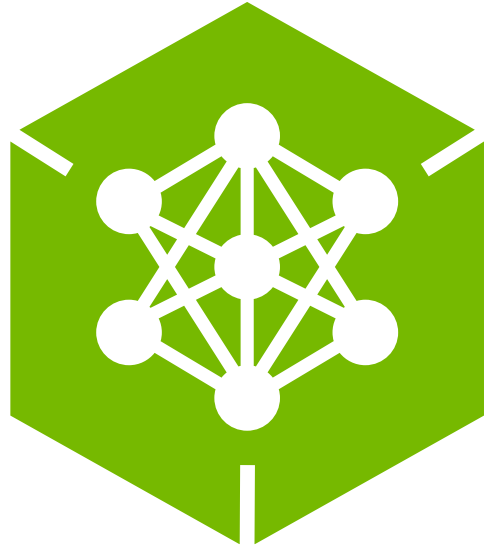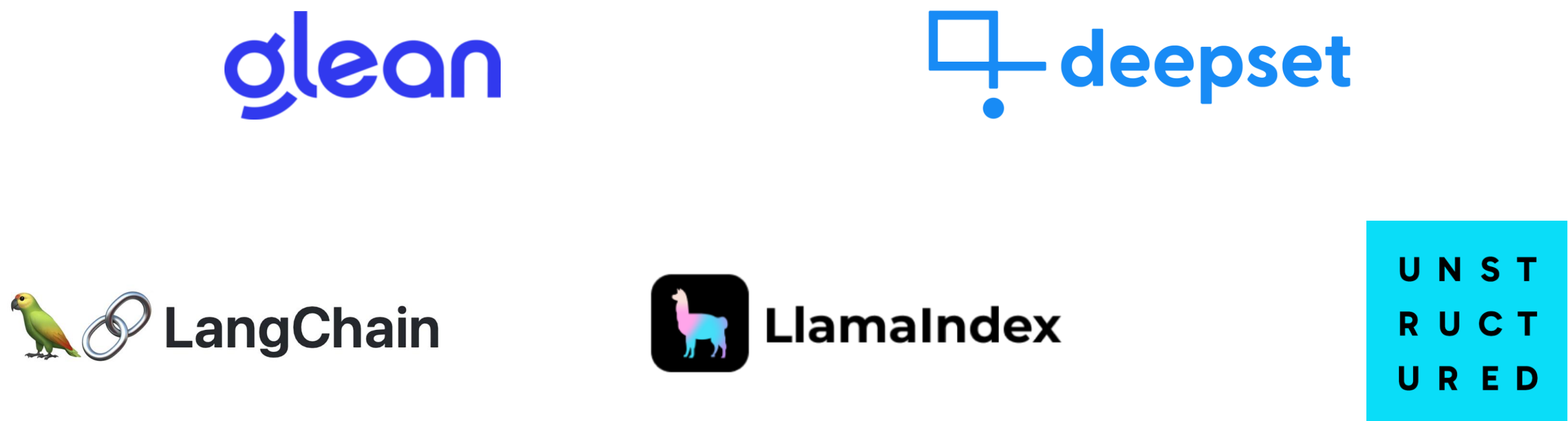Build, customize, and deploy generative AI models with NVIDIA NeMo.

# AI-Based Information Retrieval Unlocks Enterprise Knowledge

## NVIDIA NeMo Ecosystem

DATA PLATFORMS

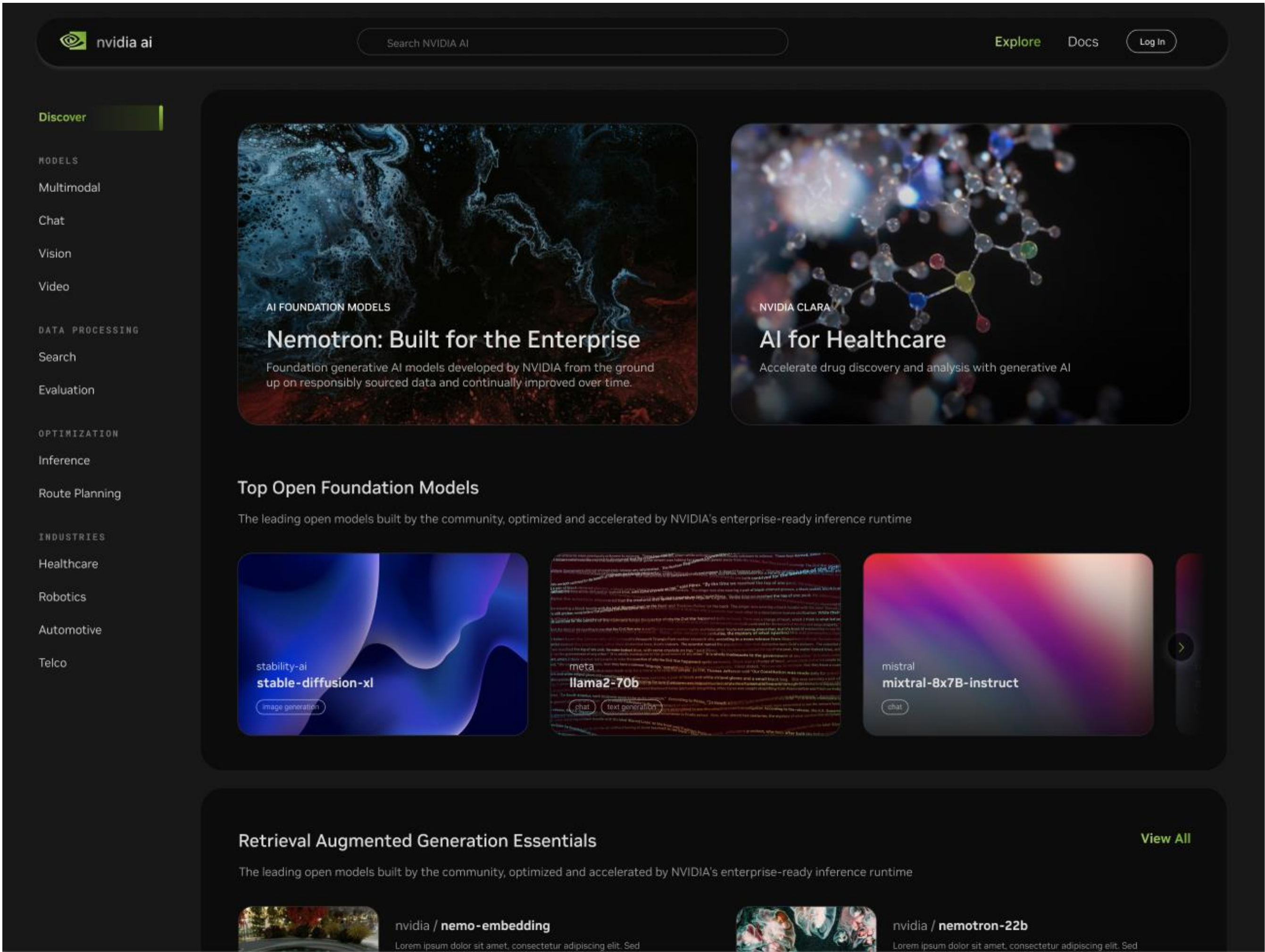RAG ECOSYSTEM

VECTOR DATABASE

# Getting Started
## Endpoints & EA

**1** NVIDIA API Catalog

*Includes all APIs across NVIDIA*



**2** Apply for Early Access

*For NeMo Retriever microservices*