# Generative AI

ICN Solutions BIM Congres – 24 September 2024

Jean-Pierre van Gastel – HP – NVIDIA tech evangelist

# AI Adoption is Growing and is Business Critical
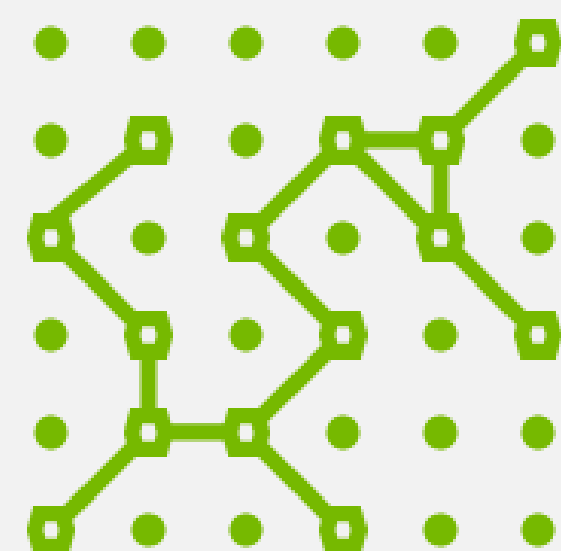## Across the hybrid cloud

### Increasing AI Adoption

**75%** of large enterprises will use AI to enhance efficiency and improve quality[1]

**56%** average adoption rate of AI by organizations globally

### Struggle with Complexity

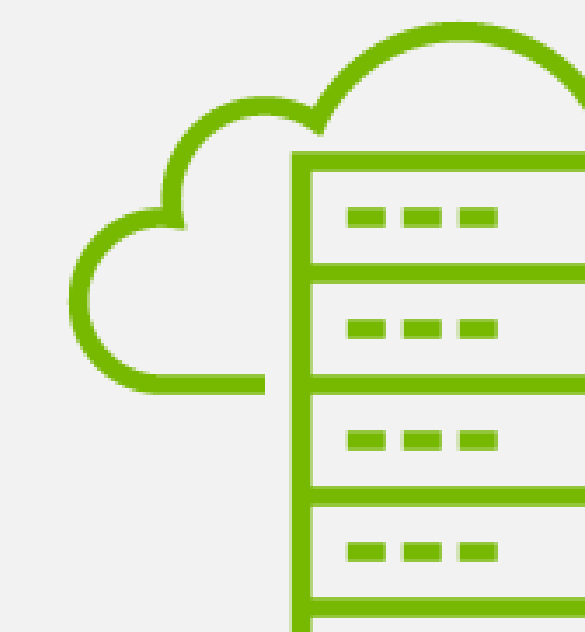**7.3** months from pilot to production[2]

**31%** have AI deployed in production

### Growing Adoption of Cloud

**90%** of enterprise cloud infrastructure will be based on public cloud providers

**50%** of all accelerated infrastructure for performance-intense computing will be cloud based

[1]IDC. "IDC FutureScape, Worldwide AI and Automation 2023 Predictions", 2022; [2]Gartner, "2023 Planning Guide for Analytics and AI", 2023; [3]Gartner, "Forecast Analysis: Cloud Infrastructure and Platform Services, Worldwide", 2021

<span>◎ nvidia.</span>

# Terminology Explained

Workload vs. workflow



## Workload

Any software program or application, that is standalone or part of a workflow, that uses compute resources to accomplish a task.

Data science, AI, and 3D graphics workloads can be accelerated by libraries and frameworks that leverage NVIDIA GPUs.

**Examples: Spark jobs, models doing video analytics, training a large language model, a text-to-speech function, video rendering**
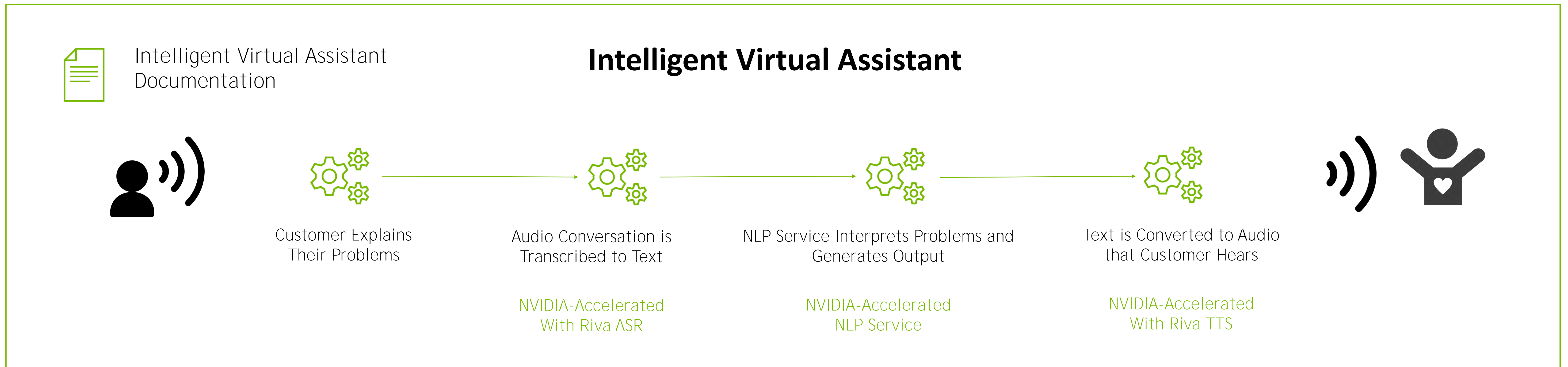
## Workflow

Multi-step process to get from initiation to completion, where each step is a unique workload. For example, the generic workflow of AI is data prep > training > simulation > inference.
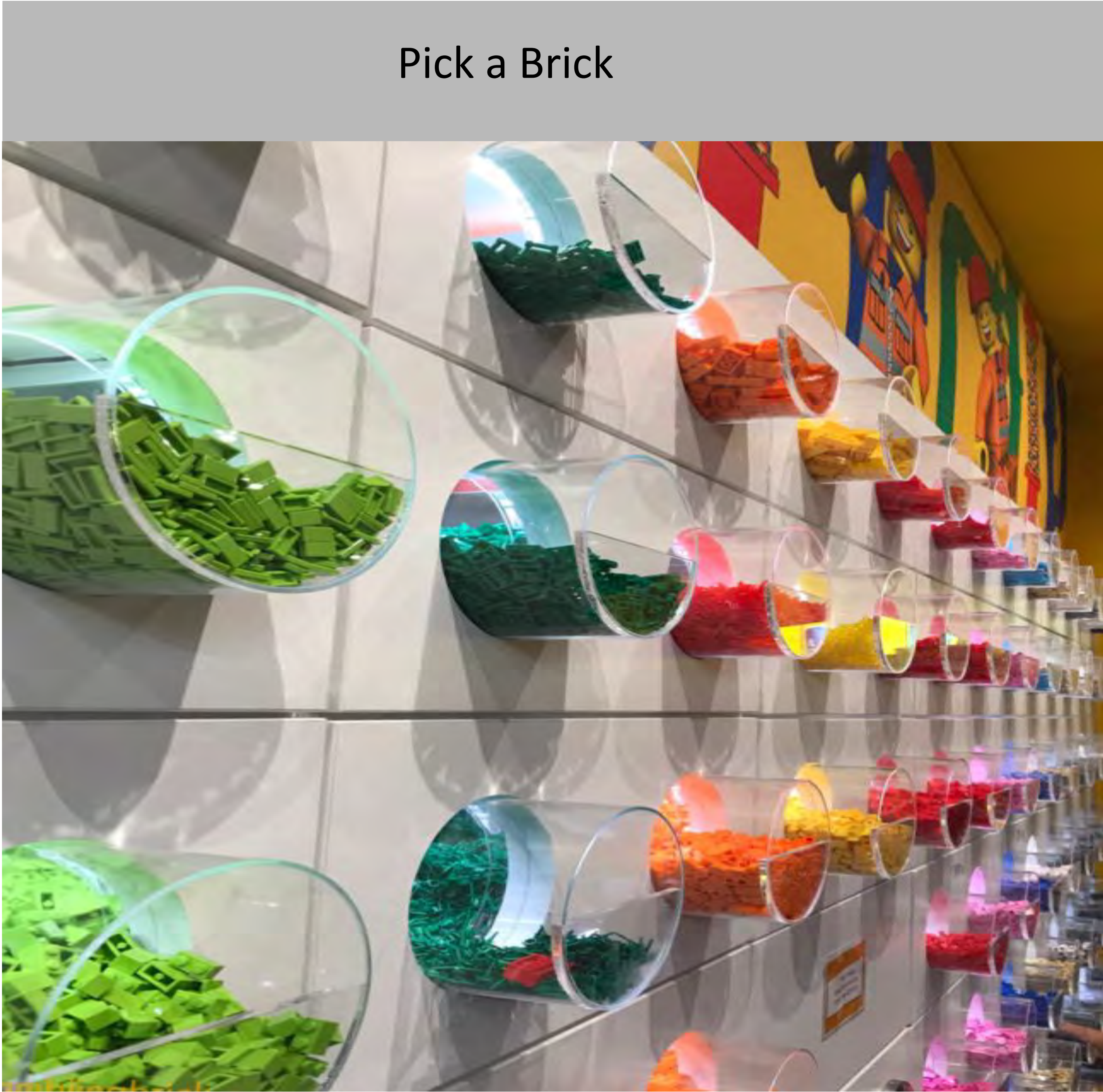
NVIDIA has assembled, tested, and documented reference workflows that can be customized by partners and customers to give them a head start at solving their specific challenge.

**Examples: Audio Transcription, Digital Fingerprinting to Detect Cybersecurity Threats, Contact Center Intelligent Virtual Assistants**
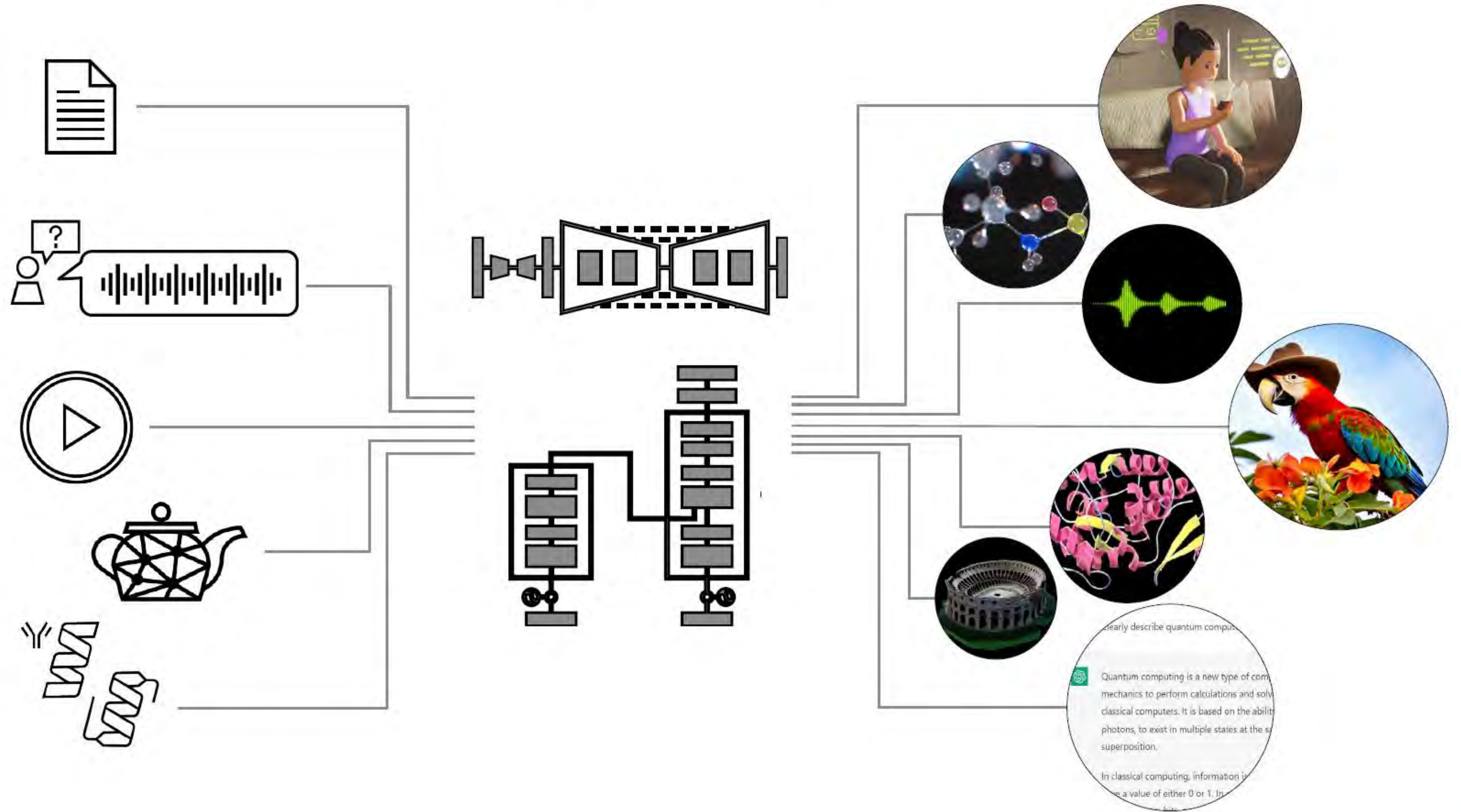
# Workflow Examples



## Making Cookies

Chocolate Chip Cookies Recipe

Gather and Prepare Ingredients → Combine Ingredients in Mixer → Scoop Dough onto Baking Sheet → Bake Cookies in the Oven → Cool Cookies on Baking Rack

## Intelligent Virtual Assistant

Intelligent Virtual Assistant Documentation

Customer Explains Their Problems → Audio Conversation is Transcribed to Text → NLP Service Interprets Problems and Generates Output → Text is Converted to Audio that Customer Hears

NVIDIA-Accelerated With Riva ASR

NVIDIA-Accelerated NLP Service

NVIDIA-Accelerated With Riva TTS

# Two Ways to Build with Legos



Pick a Brick

DIY, unlimited options



Lego Kit

Components you need to build the AI solution

# Two Ways to Build with NVIDIA AI



## Public Catalog on NGC

DIY, unlimited options

## NVIDIA AI Workflows

Intelligent Virtual Assistant

Digital Fingerprinting Cybersecurity Threat Detection

Audio Transcription

Product Recommendations
xxx

Components you need to build the AI application

# What is Generative AI?

# Generative AI
## From Research To Production In 5 Years



ChatGPT, LLAMA2

Stable Diffusion etc

Adobe Photoshop Generative Fill

# Generative AI is Transforming Every Industry



| 3D VFX & Game Design | Architecture & Interior Design | Fashion & Product Design | Photography & Photo Editing | Marketing and Advertising | Manufacturing |

Generate textures and backgrounds

Create floorplans and explore architectural styles

Inspire unique design concepts

Background and object replacement

Create elements & reusable motifs

Design parts
Explore structures & solutions

# Enterprise are on the Generative AI Journey

**2022**

**2023**

**2024**

## 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.
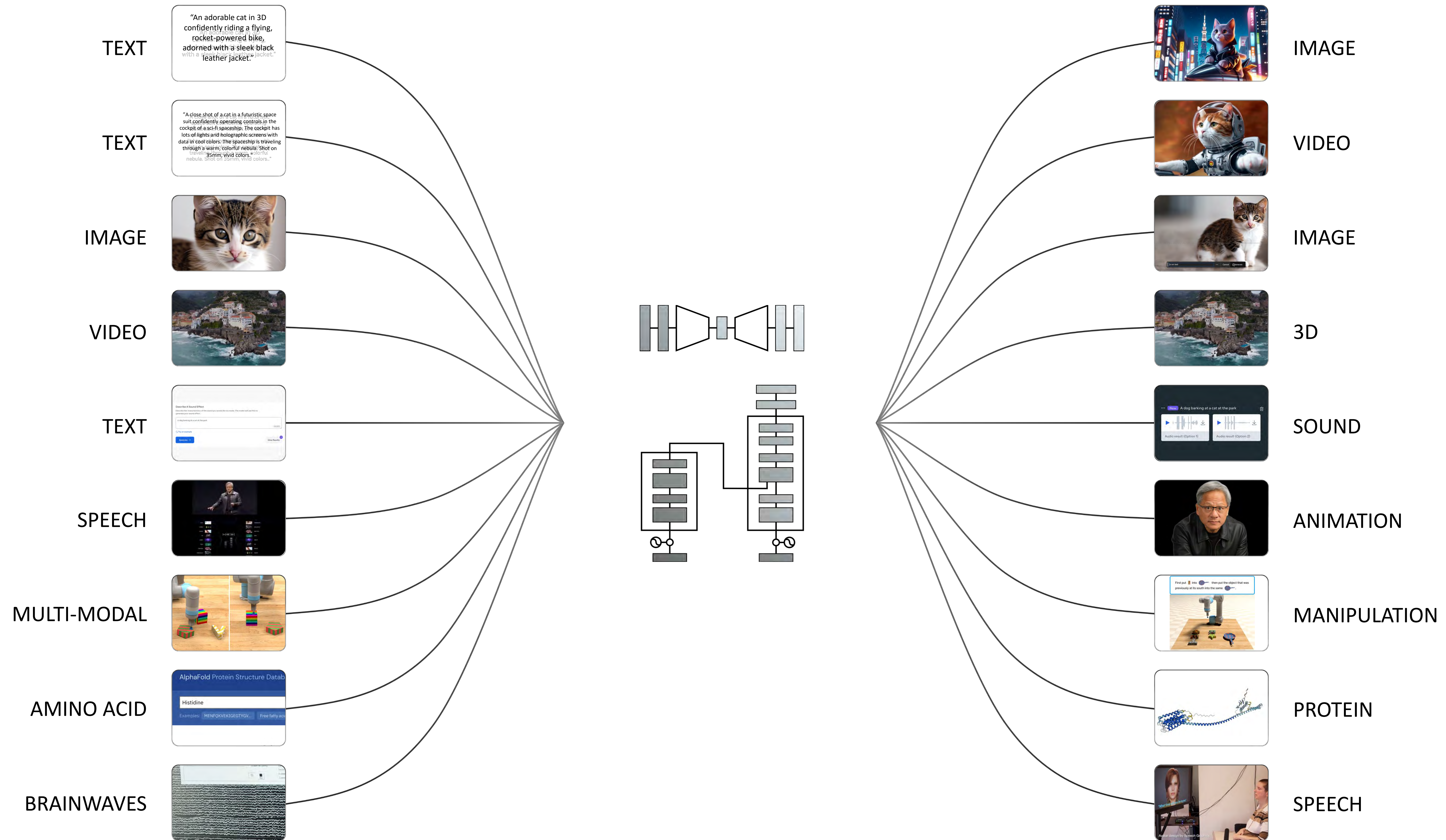
## Experimentation

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

## Production

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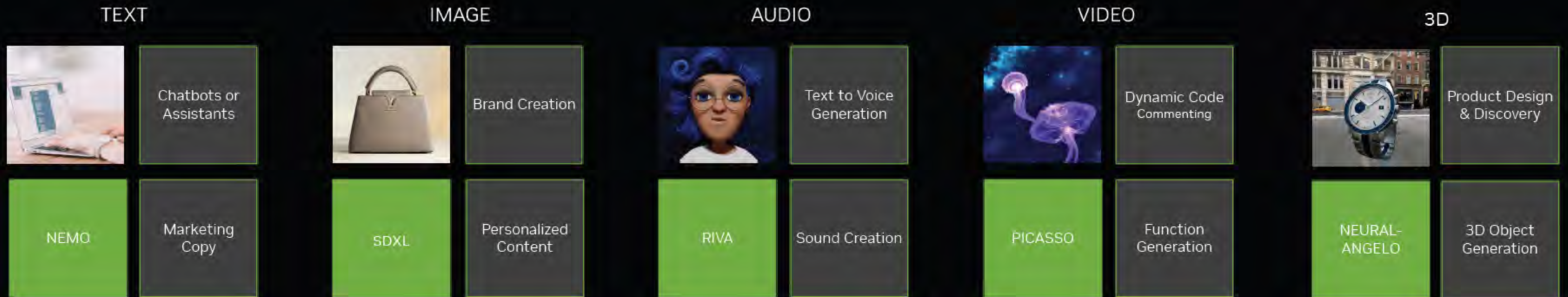
**NVIDIA**

# Generative AI Can Learn and Understand Everything
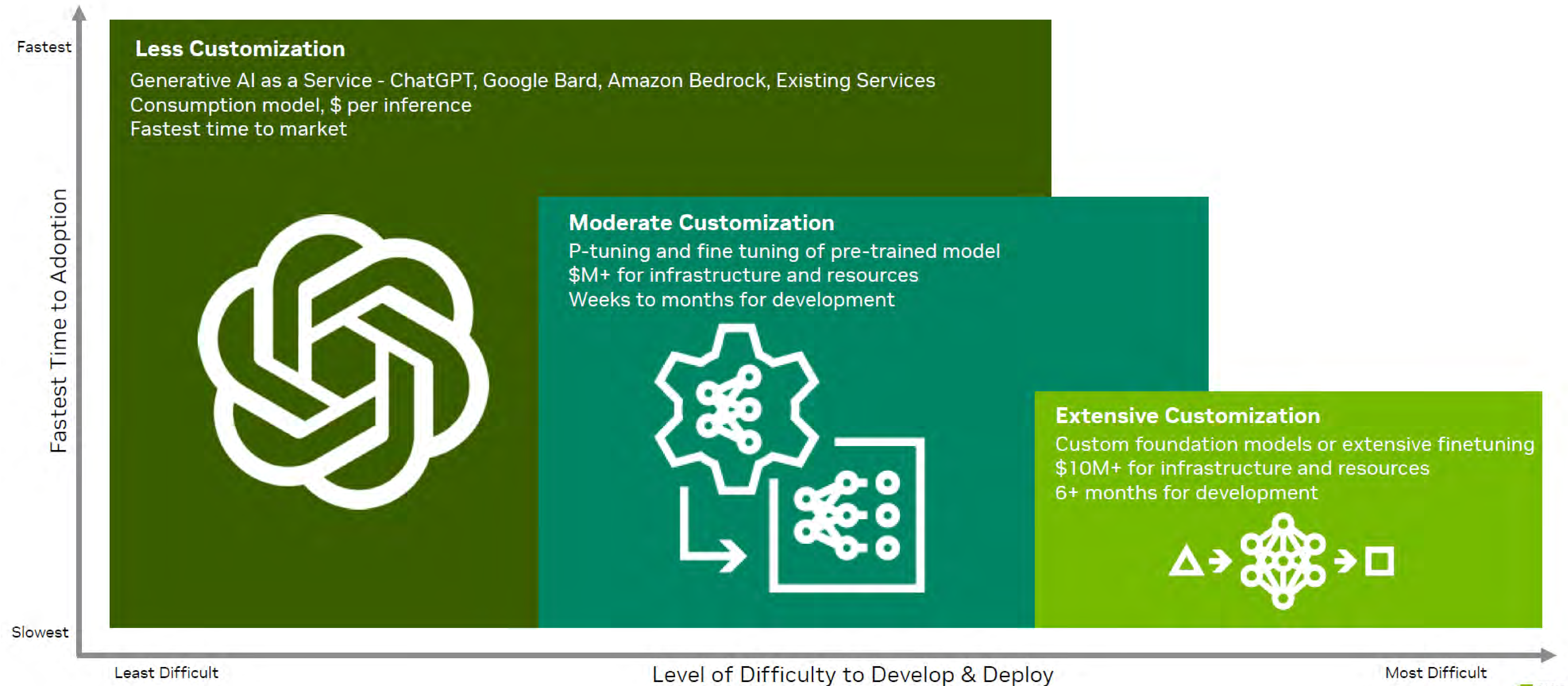
# Autodesk 3ds Max – tyFlow demo

# Multimodal Generative AI is Transforming Industries



Enterprises that adopt next-generation AI like LLMs and Generative AI are **2.6X more likely to increase revenue by 10% or more** but must invest in their AI infrastructure to fully reap the benefits.

-Accenture Research. Breakthrough Innovation: Is your organization equipped for breakthrough innovation? WEF 2023.

# How Enterprises are Using Generative AI

Fastest ↑

Fastest Time to Adoption

Slowest

**Less Customization**

Generative AI as a Service - ChatGPT, Google Bard, Amazon Bedrock, Existing Services
Consumption model, $ per inference
Fastest time to market

**Moderate Customization**

P-tuning and fine tuning of pre-trained model
$M+ for infrastructure and resources
Weeks to months for development

**Extensive Customization**

Custom foundation models or extensive finetuning
$10M+ for infrastructure and resources
6+ months for development

Least Difficult

**Level of Difficulty to Develop & Deploy**

Most Difficult

nVIDIA.

# Enterprise Generative AI Use Cases Require Domain Specific Knowledge

## Foundation Model Response

*"When did I last send a payment to my credit card company?"*

→

*"I was trained 2 months ago and do not have the current data "*

## Custom Model Response

*"When did I last send a payment to my credit card company?"*

→

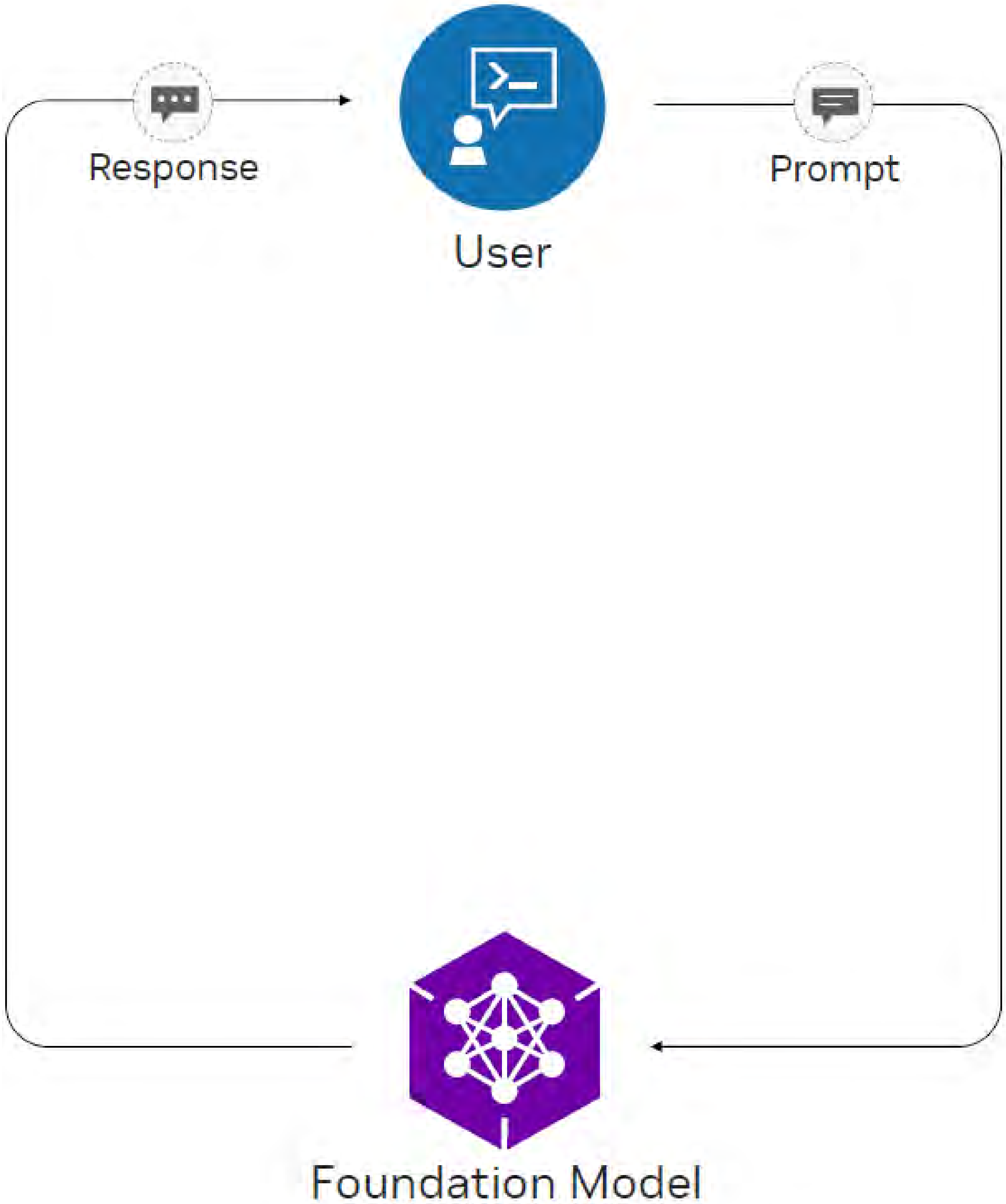*"The last payment was sent on May 27, 2023."*

Dataset

## 70%

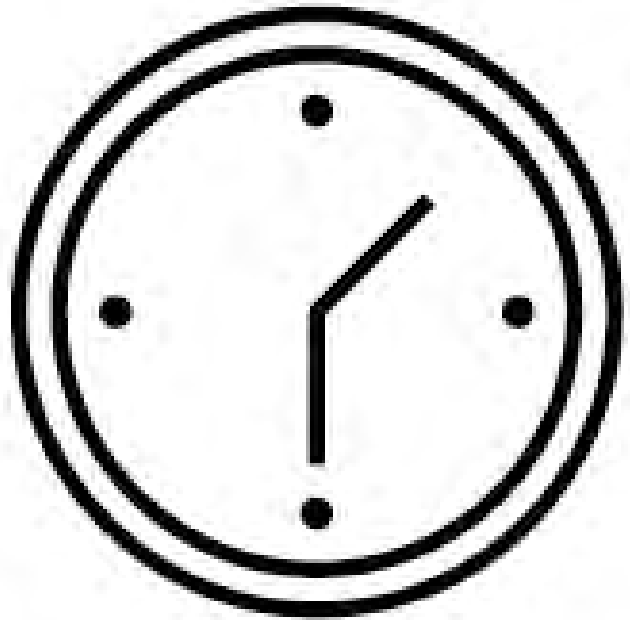**Of enterprise data is untapped**
Unlock many new opportunities for greater intelligence

**Less frequent re-training**
Significant cost and time savings in long-run to maintain LLMs

# LLMs are Powerful Tools but Not Accurate Enough for Enterprise

## Without a connection to enterprise data sources, LLMs cannot provide accurate information

Response

User

Prompt

Foundation Model
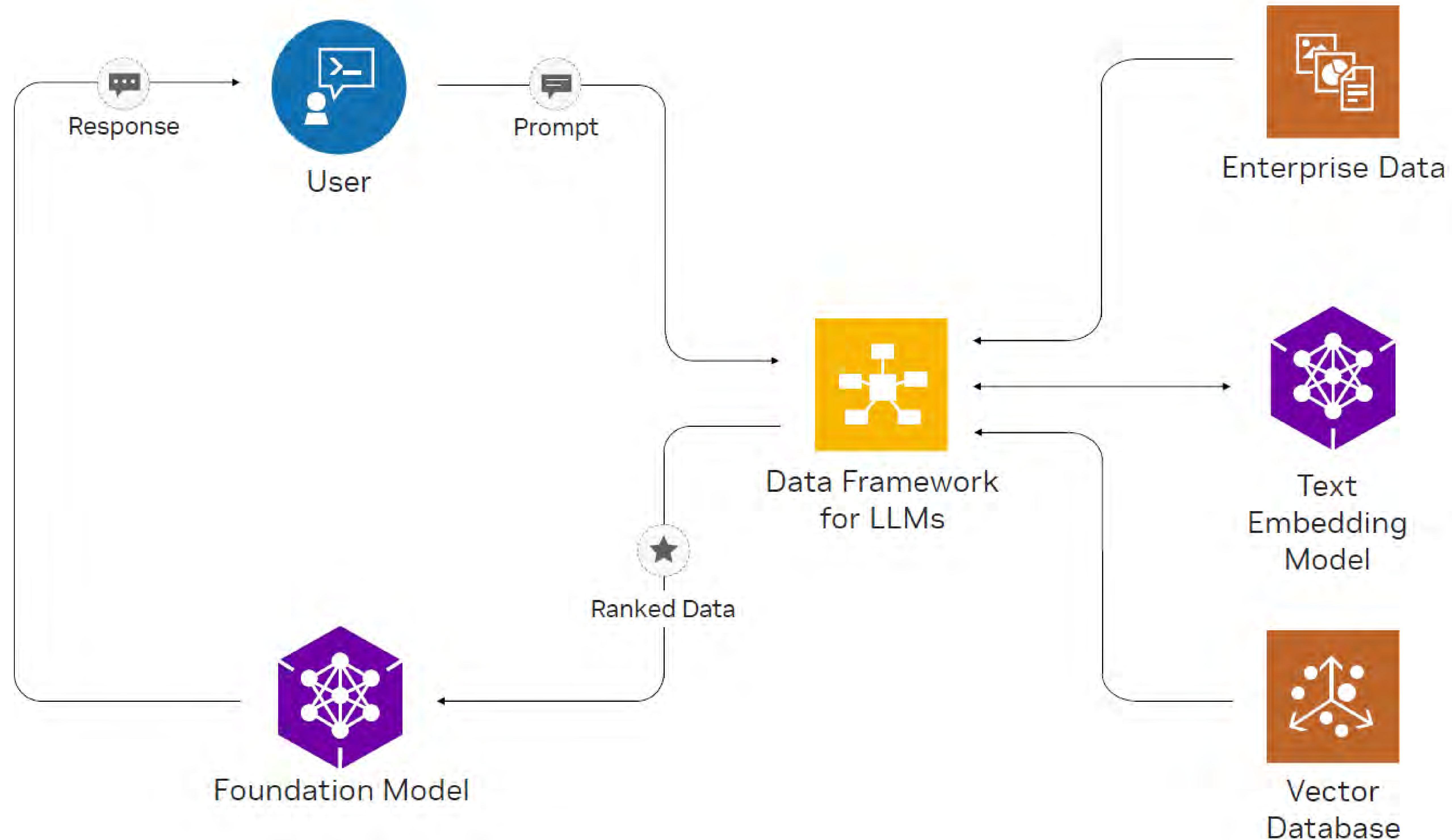
Lacking proprietary knowledge

Risk of outdated information

Hallucinations

NVIDIA

# Retrieval Augmented Generation Lets Enterprises Talk to Their Data

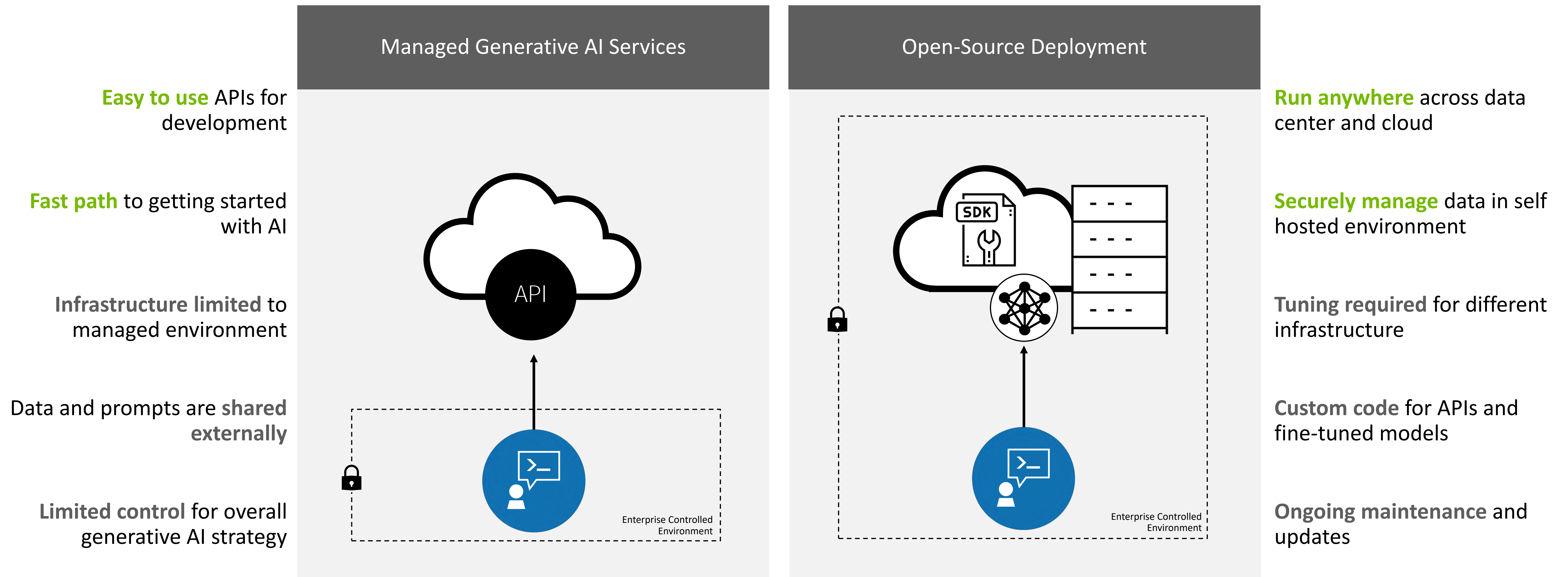## Enable LLMs to provide up to date and domain specific answers

# NVIDIA ChatRTX

Free technology RAG demo that runs on your workstation

https://www.nvidia.com/en-us/ai-on-rtx/chatrtx/

# Enterprises Face Challenges Experimenting with Generative AI

## Organizations must choose between ease of use and control



**Managed Generative AI Services**

**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

Enterprise Controlled Environment

**Open-Source Deployment**

**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

Enterprise Controlled Environment

# NVIDIA AI Foundation Models and Endpoints

Fast-track custom generative AI models for enterprise applications
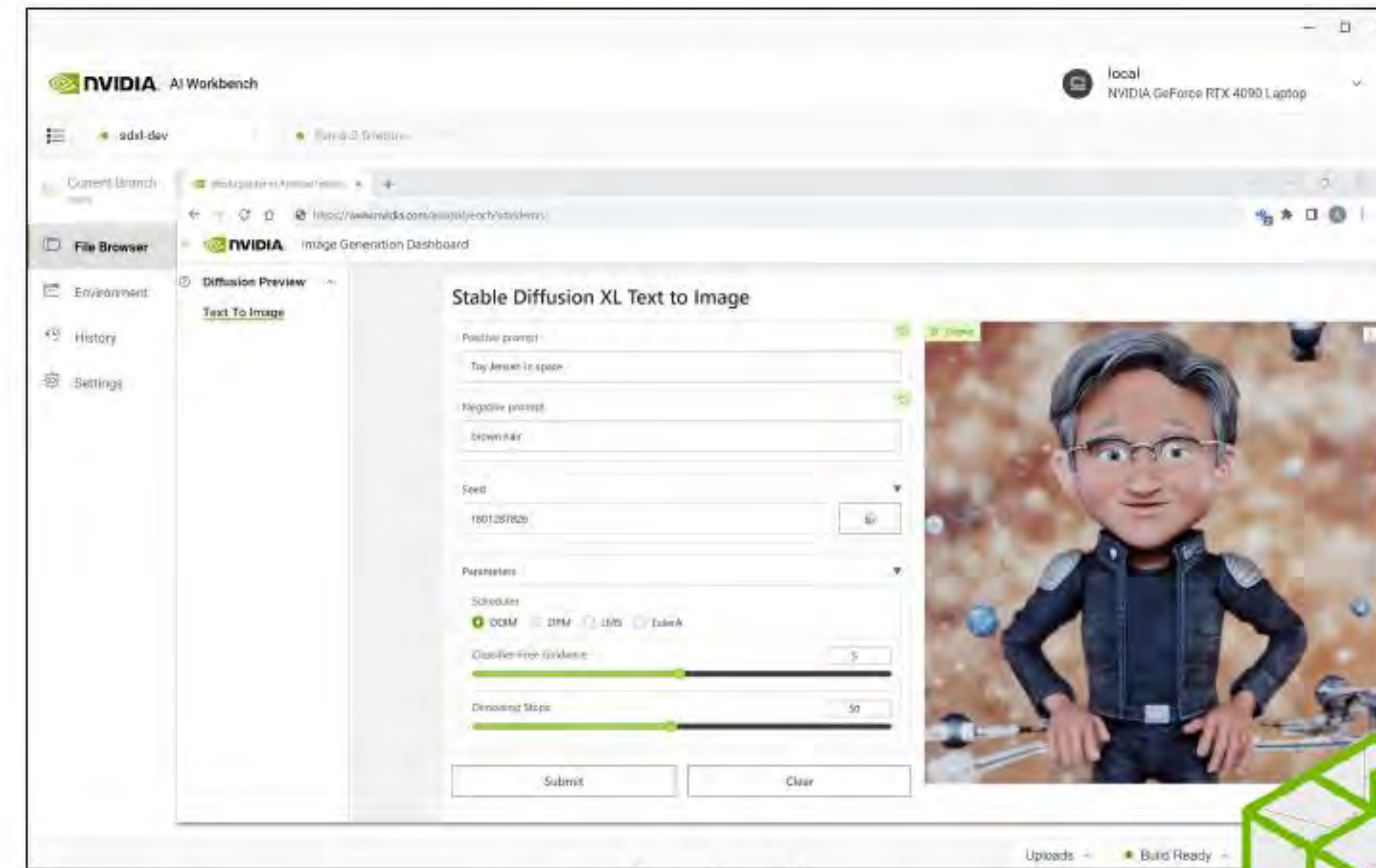


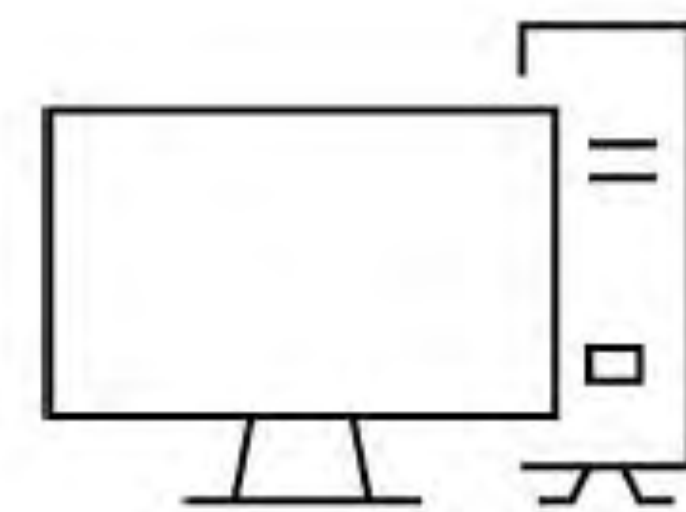Enterprise-ready, performance optimized models from NVIDIA and the community

Experience foundation models running on the NVIDIA AI stack via API endpoints

# NVIDIA AI Workbench

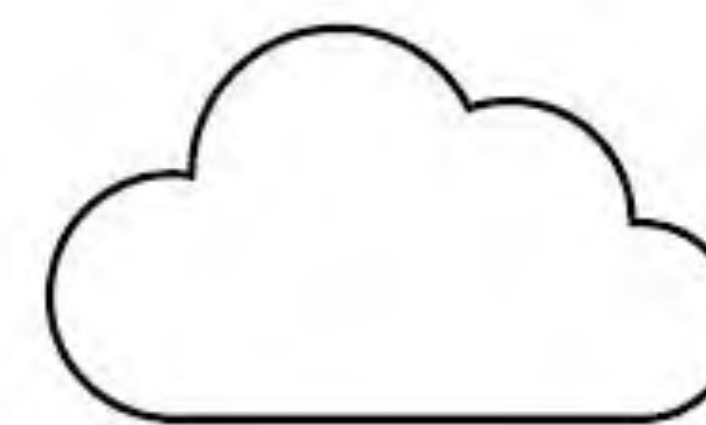## Enables anyone with access to a GPU to be a generative AI creator



- Create projects for tuning and deployment of generative AI and LLMs

- Move projects between PCs and workstations, data centers, public clouds, and NVIDIA DGX Cloud

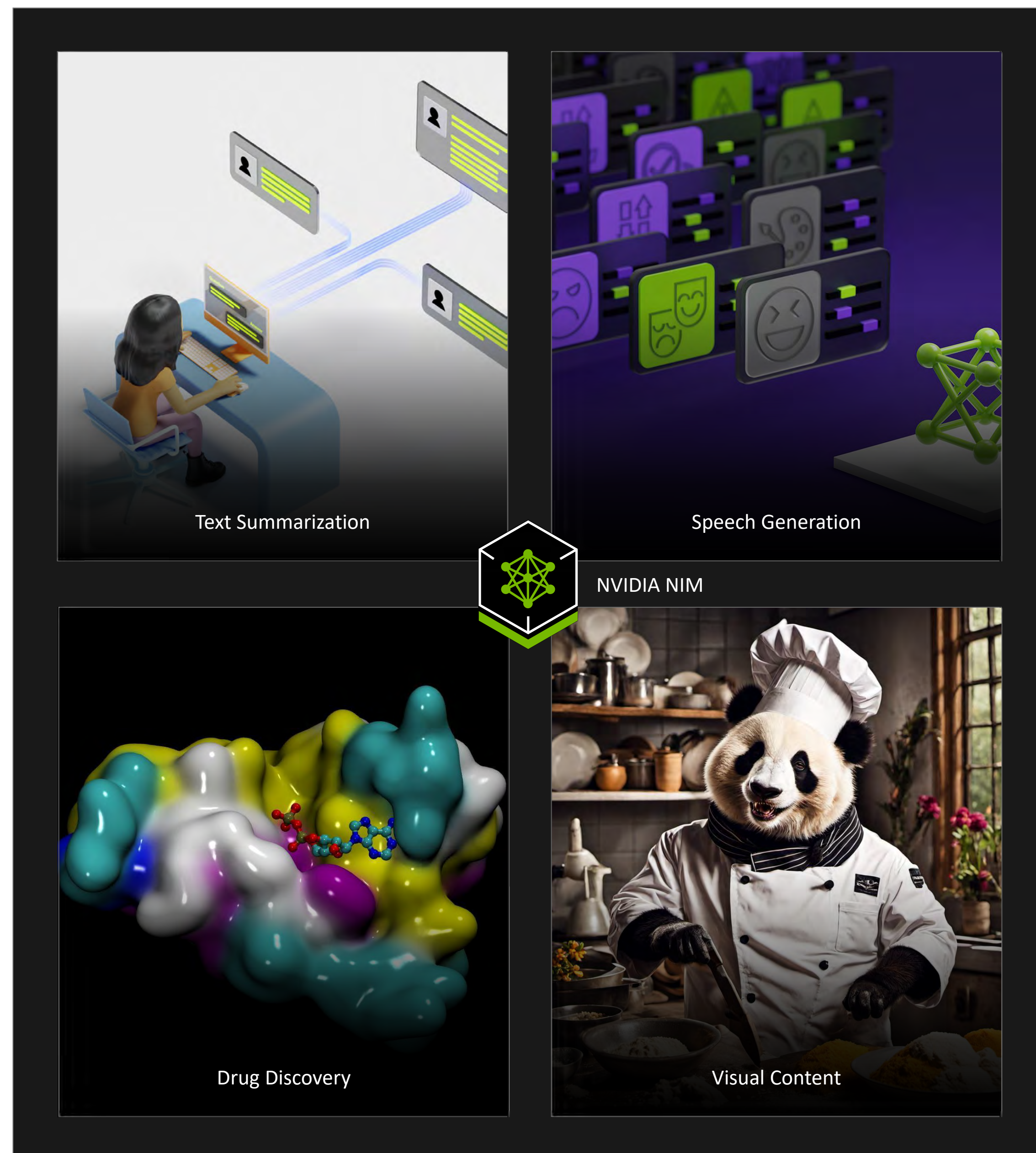- Easily start with pre-built project examples
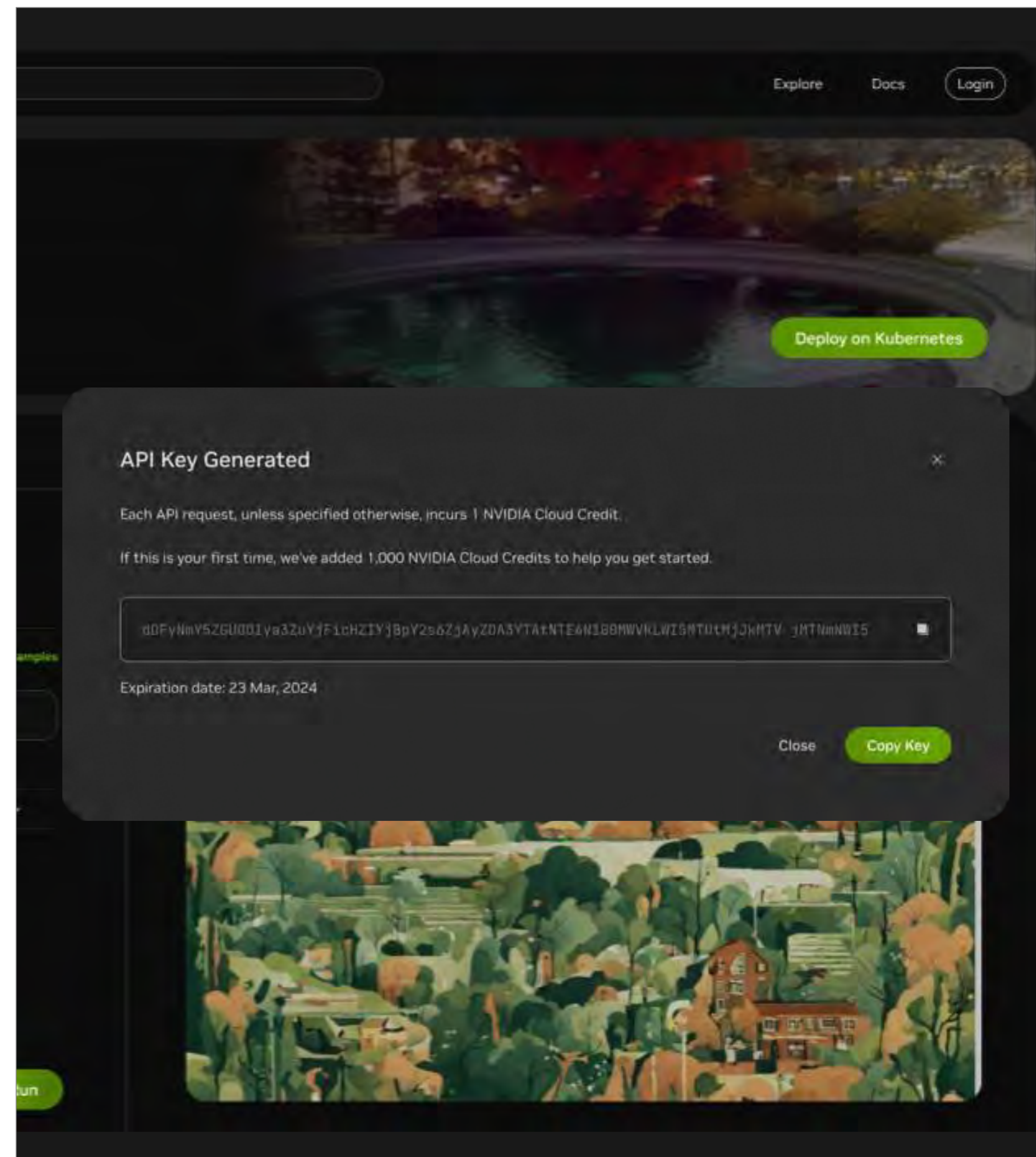
**PCs & WORKSTATIONS**     **DATA CENTERS**     **CLOUDS**
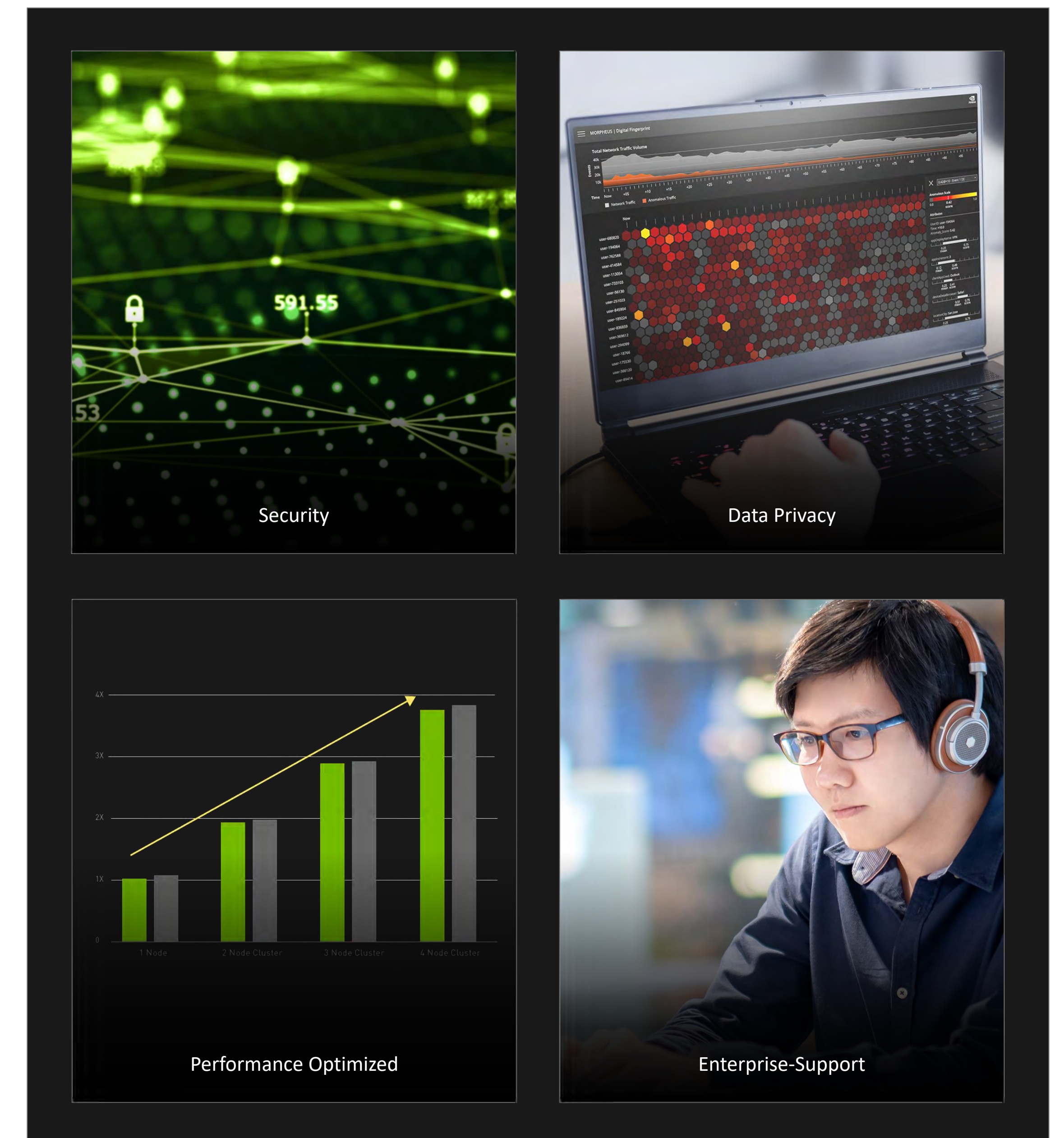
# Experience and Run Enterprise Generative AI Models Anywhere

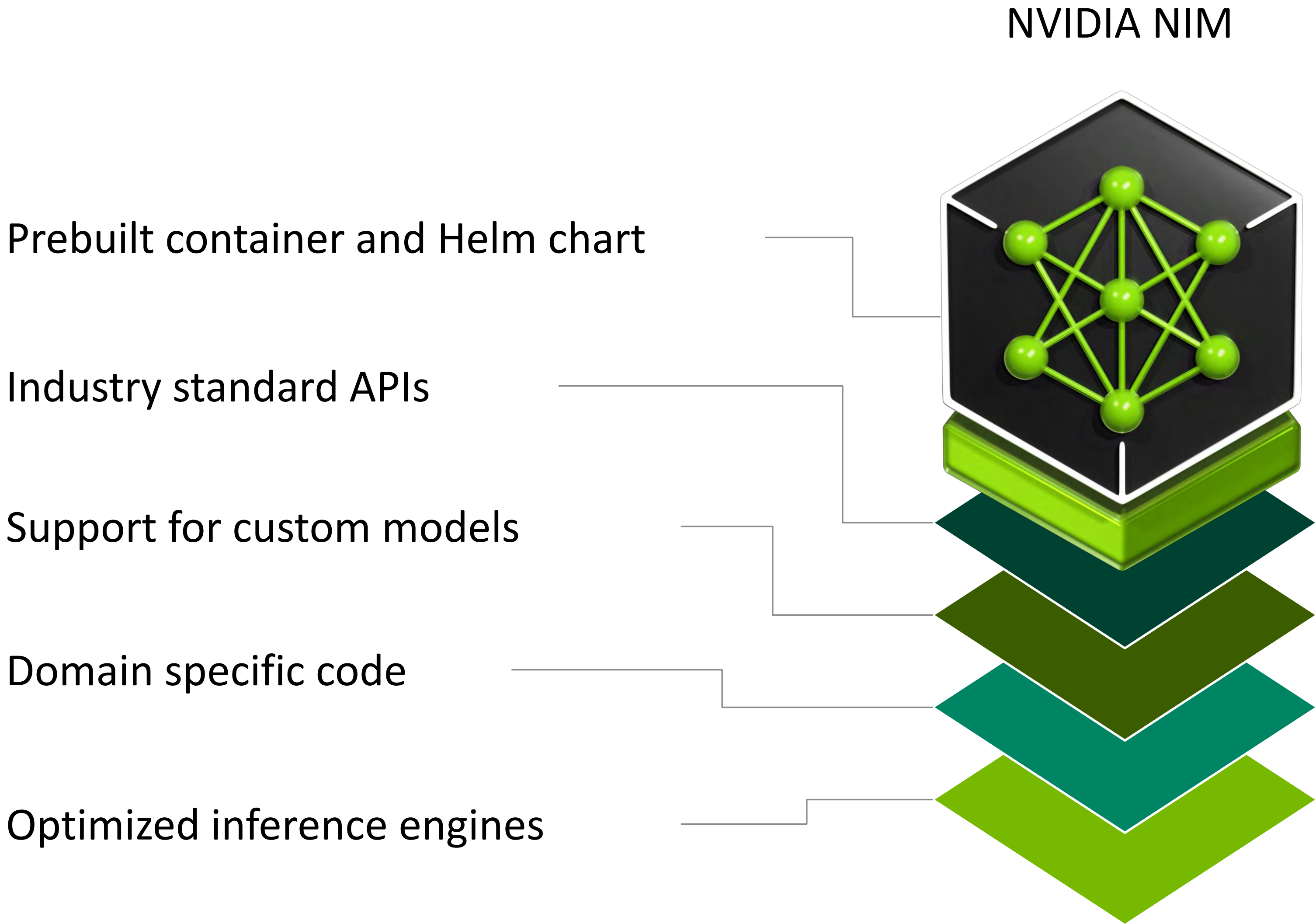## Use NVIDIA API catalog to get access to NVIDIA NIM



Text Summarization

Speech Generation

NVIDIA NIM

Drug Discovery

Visual Content

**Experience Models**



Deploy on Kubernetes

**API Key Generated**

Each API request, unless specified otherwise, incurs 1 NVIDIA Cloud Credit.

If this is your first time, we've added 1,000 NVIDIA Cloud Credits to help you get started.

mDFyNmVSZGUODIyw3ZuYjFichZIYjBpY2s6ZjAyZDA3YTAtNTE6N1B0HWVKLWIS0TD1tMjJxHTV-jHTNmNBI5

Expiration date: 23 Mar, 2024

Close    Copy Key

**Prototype with APIs**



Security

Data Privacy

Performance Optimized

Enterprise-Support

**Deploy with NIMs**

NVIDIA

# NVIDIA NIM Optimized Inference Microservices

## Accelerated runtime for generative AI

NVIDIA NIM

Prebuilt container and Helm chart

Industry standard APIs

Support for custom models

Domain specific code

Optimized inference engines

**Deploy anywhere and maintain control** of generative AI applications and data

**Simplified development** of AI application that can run in enterprise environments

**Day 0 support** for all generative AI models providing choice across the ecosystem

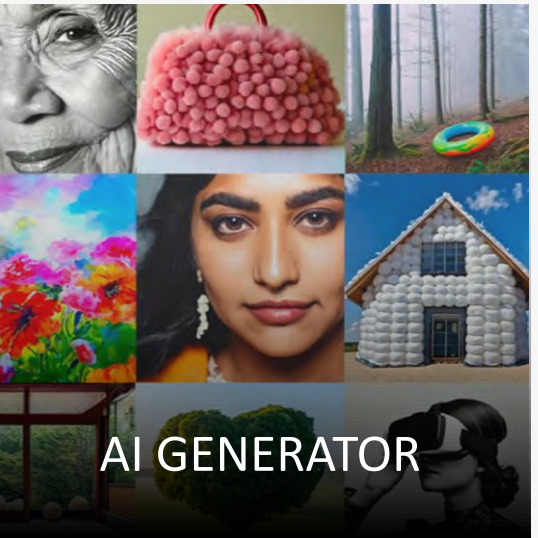**Improved TCO** with best latency and throughput running on accelerated infrastructure

**Best accuracy** for enterprise by enabling tuning with proprietary data sources

**Enterprise software** with feature branches, validation and support

# Inference Microservices for Generative AI

NVIDIA NIM is the fastest way to deploy AI models on 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-Certified Systems through leading partners

# NVIDIA NIM for Every Domain

## LANGUAGE NIMs

| | | |
|---|---|---|
| Code Llama 70B | Cohere 35B | Gemma 7B |
| Jamba | Llama 2 70B | Mistral 7B |
| Mixtral 8x7B | Nemotron-3 22B Persona | Phi-2 |

## VISUAL / MULTIMODAL NIMs

| | | |
|---|---|---|
| Adept 110B | Deplot | Edify. Getty |
| Edify. Shutterstock | FuYu 8B, 55B | Kosmos-2 |
| NeVA | SDXL 1.0 | SDXL Turbo |

## DIGITAL HUMAN NIMs

| | |
|---|---|
| Audio2Face | Riva ASR |

## OPTIMIZATION / SIMULATION NIMs

| | |
|---|---|
| cuOpt | Earth-2 |

## DIGITAL BIOLOGY NIMs

| | |
|---|---|
| DeepVariant | |
| DiffDock | ESMFold |
| MolMIM | Vista 3D |

## APPLICATION NIMs

| | | |
|---|---|---|
| Llama Guard | Retrieval Embedding | Retrieval Reranking |

NVIDIA

# HP Z Captis

## Capture All
## the World's Materials

Digitize materials from anywhere in minutes,[1] in a portable device equipped with a polarized and photometric computer vision system for efficiency and accuracy.

Learn More