NVIDIA

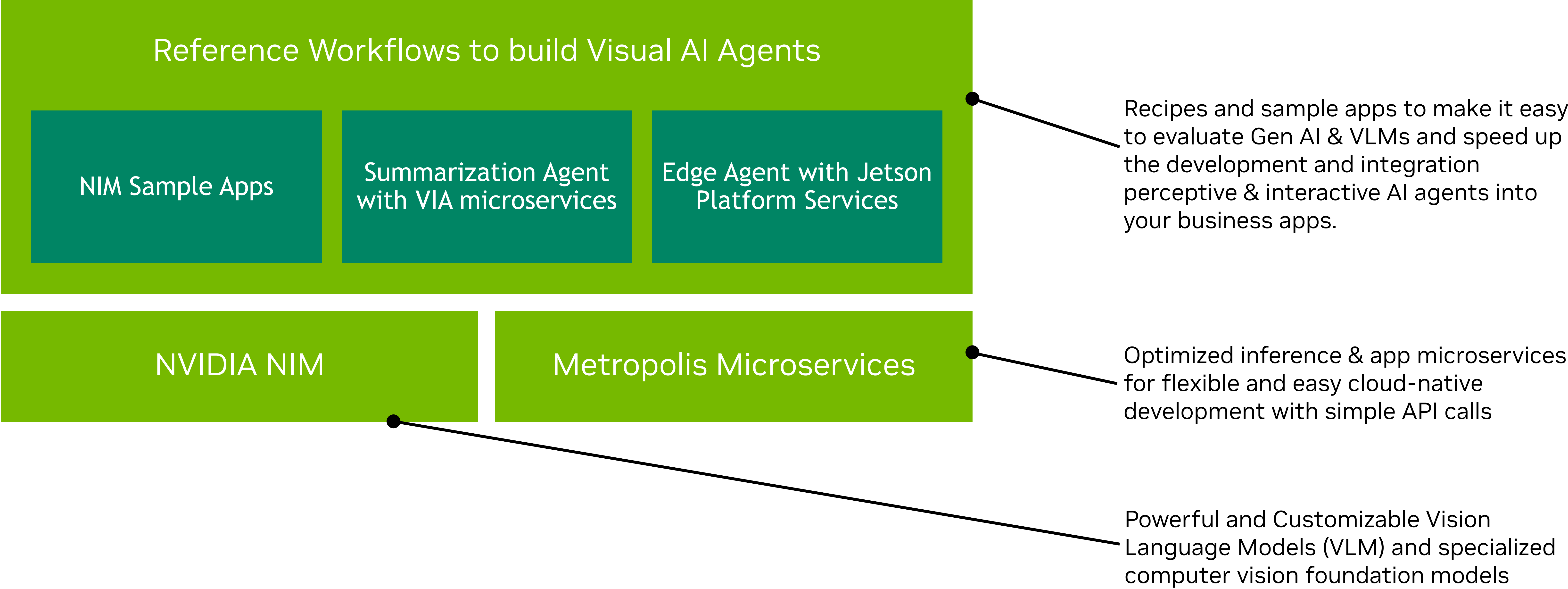# Building Visual AI Agents with Generative AI and NVIDIA NIM

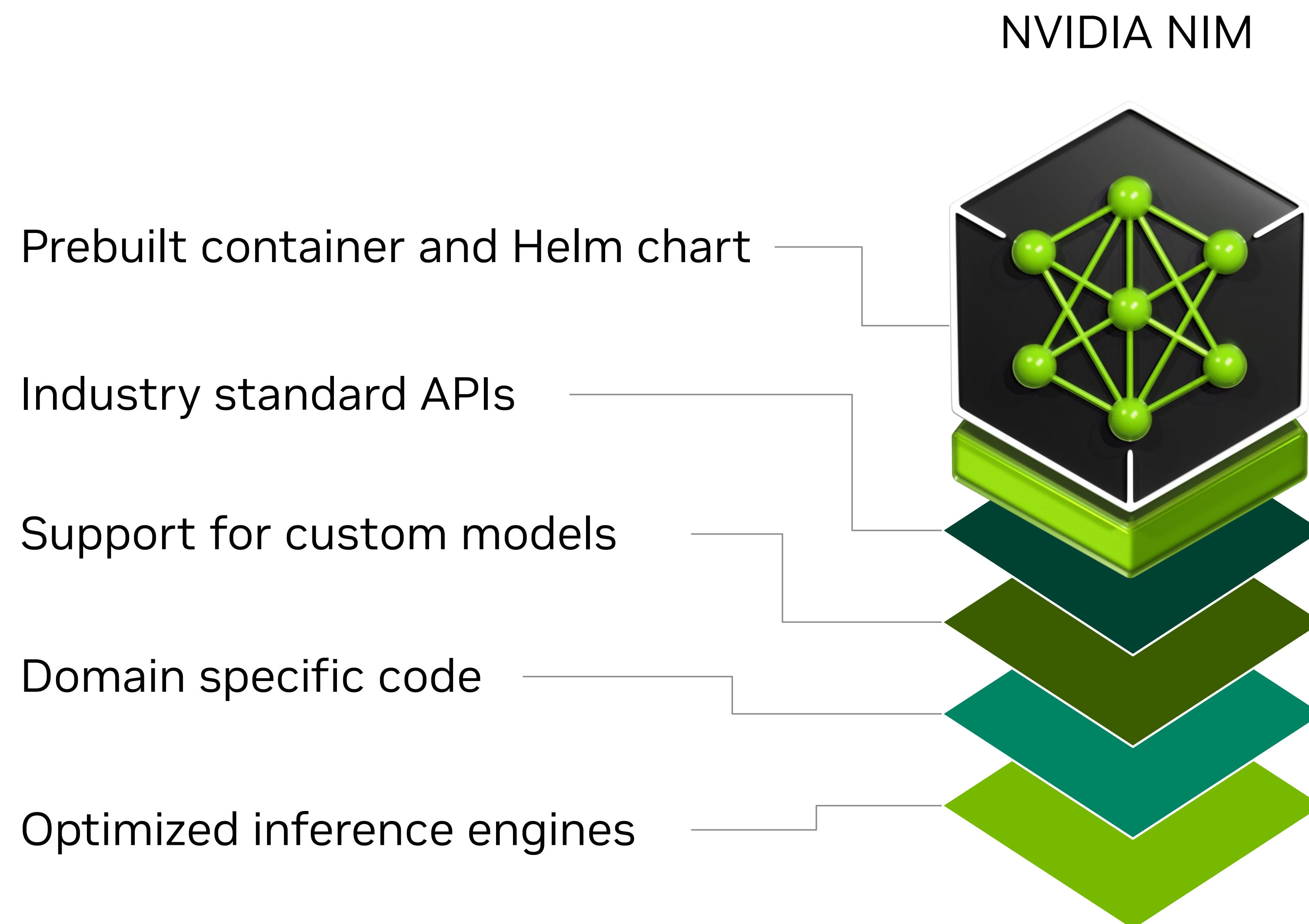NVIDIA Metropolis

# NVIDIA Metropolis Generative AI Stack



**Reference Workflows to build Visual AI Agents**

NIM Sample Apps

Summarization Agent with VIA microservices

Edge Agent with Jetson Platform Services

Recipes and sample apps to make it easy to evaluate Gen AI & VLMs and speed up the development and integration perceptive & interactive AI agents into your business apps.

**NVIDIA NIM**

**Metropolis Microservices**

Optimized inference & app microservices for flexible and easy cloud-native development with simple API calls

Powerful and Customizable Vision Language Models (VLM) and specialized computer vision foundation models

# 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 applications that can run in enterprise environments

**Day 0 support** for all publicly available models providing choice across the ecosystem

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

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

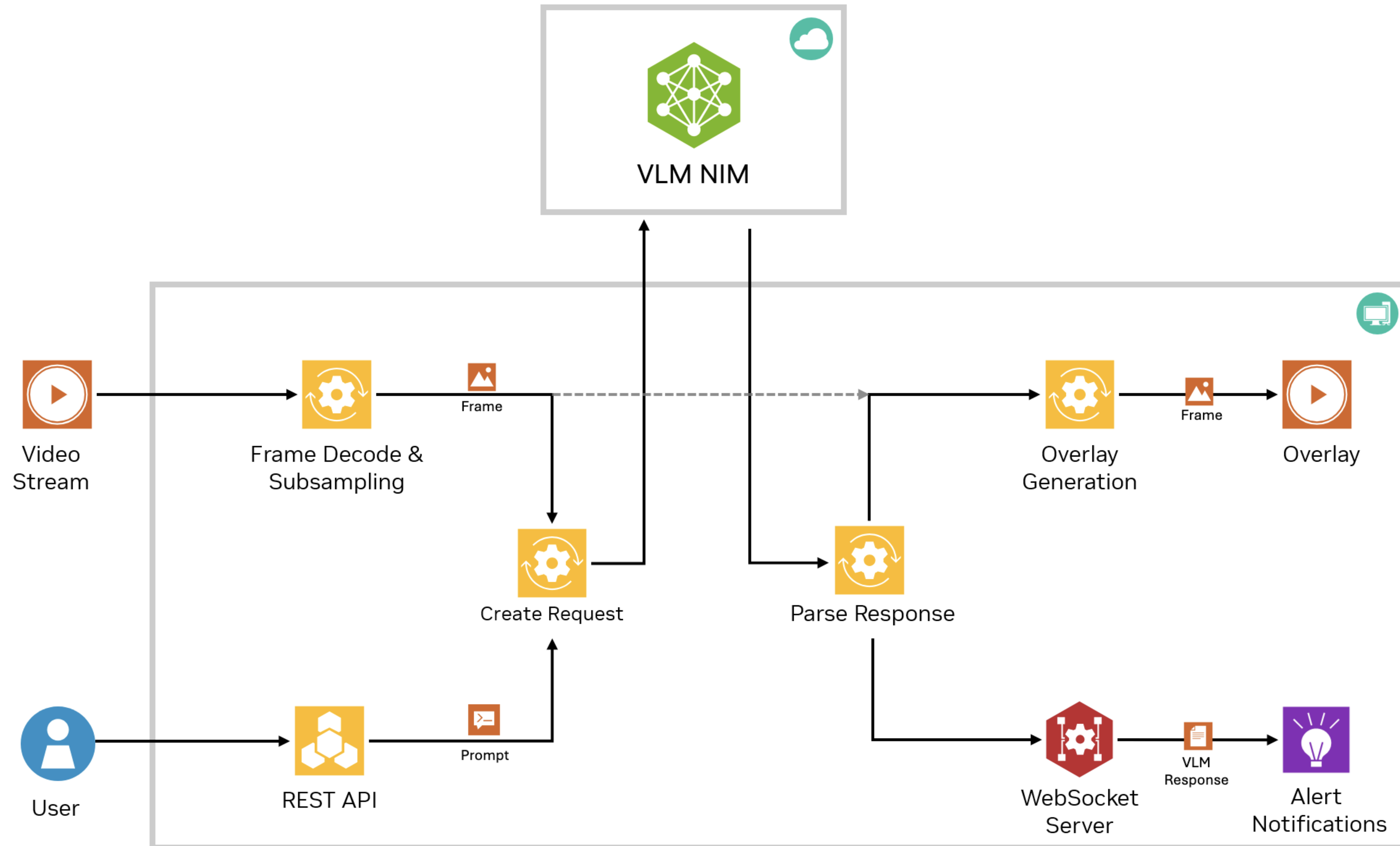**Enterprise software** with feature branches, validation and support

Microsoft Azure

aws

Google Cloud

ORACLE®

DGX & DGX Cloud

DELL Technologies

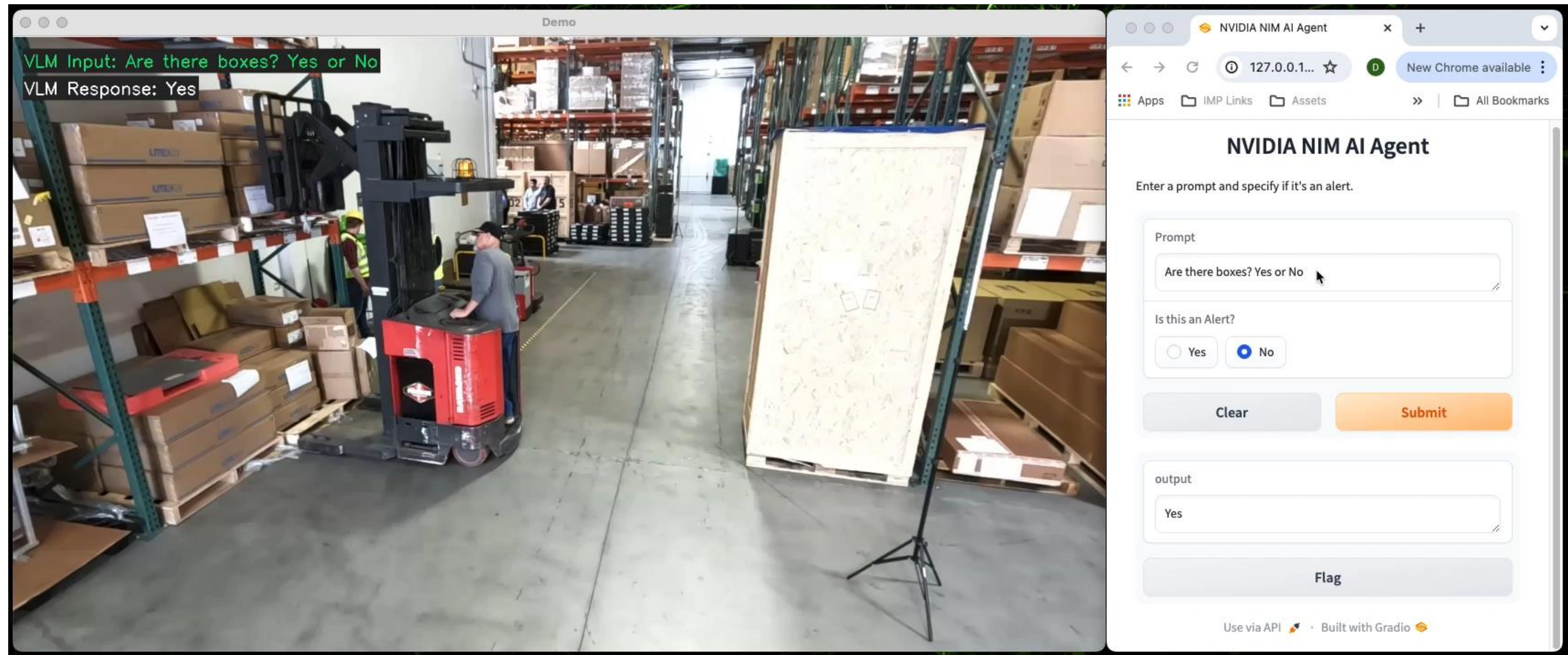Hewlett Packard Enterprise

Lenovo

SUPERMICRO

NVIDIA

# Metropolis Reference Workflows & Sample Apps

Recipes to Build a Wide Range of AI Agents

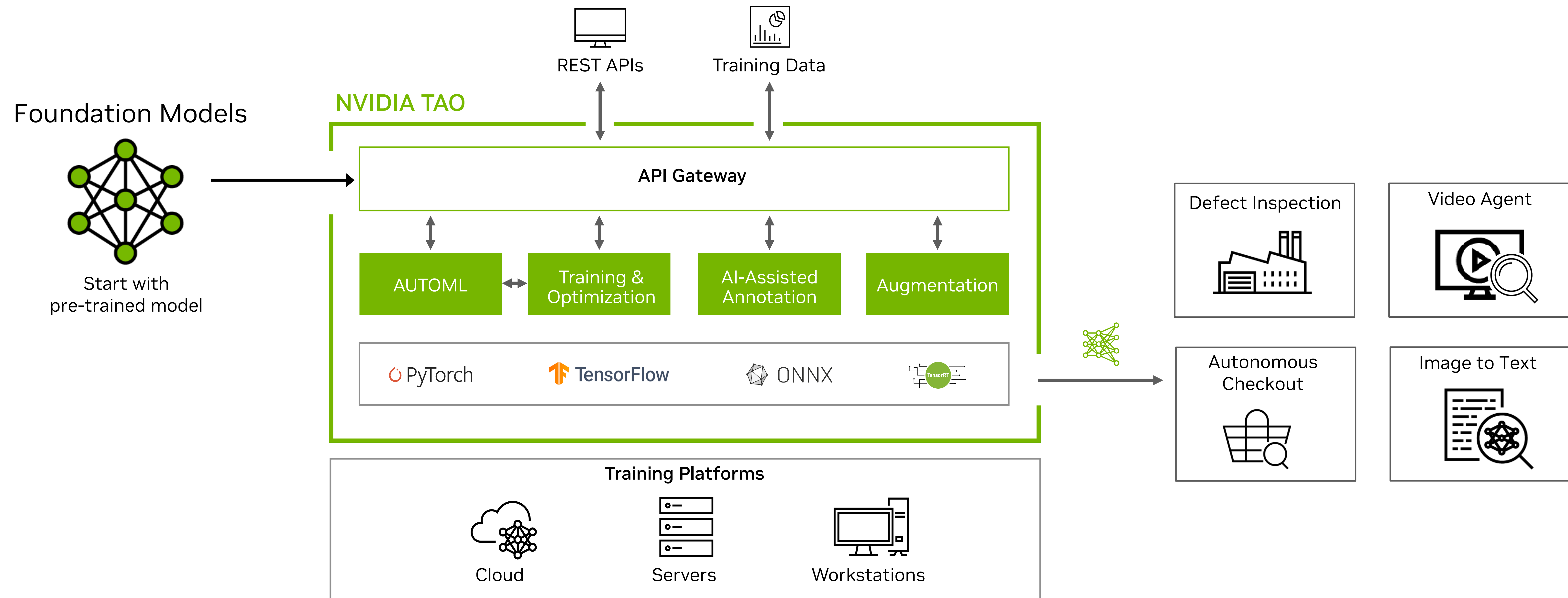# Evaluate and Experiment with Vision Language Models (VLM)

## From simple Jupyter Notebooks to Nearly Production-ready Sample apps



https://github.com/NVIDIA/metropolis-nim-workflows

# Get Access to and Customize / Tune Models with NVIDIA TAO

## State of the Art Models, Training and Customization for Vision AI

REST APIs          Training Data

**NVIDIA TAO**

Foundation Models

Start with
pre-trained model

### API Gateway

| AUTOML | Training & Optimization | AI-Assisted Annotation | Augmentation |

PyTorch          TensorFlow          ONNX          TensorRT

**Training Platforms**

Cloud          Servers          Workstations

Defect Inspection

Video Agent

Autonomous Checkout

Image to Text

**Full fine-Tuning**
Update weights of entire model
including the Foundation backbone

**Last layer or Head fine-tuning**
Freeze the Foundation backbone
and fine-tune the last few layers

**In-context Learning**
Use visual prompting and model chaining
to improve contextual awareness

NVIDIA.