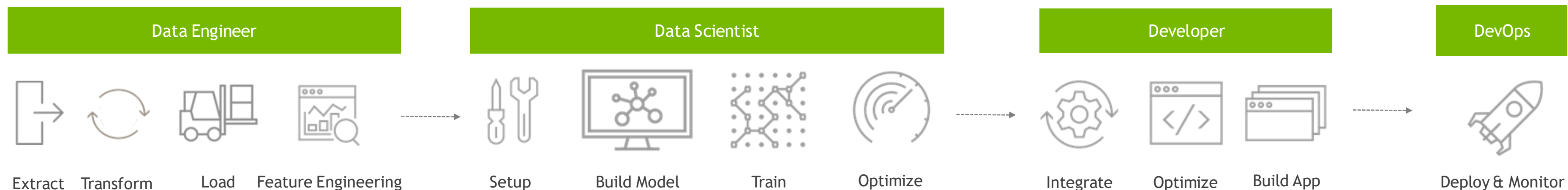# BUILD AI SOLUTIONS WITH NVIDIA NGC AND RED HAT OPENSHIFT
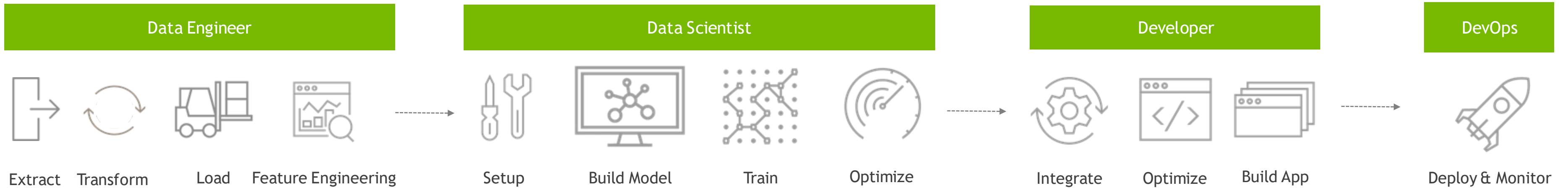
Abhishek Sawarkar, Chintan Patel. NVIDIA
Deepthi Dharwar, Diane Feddema. Red Hat

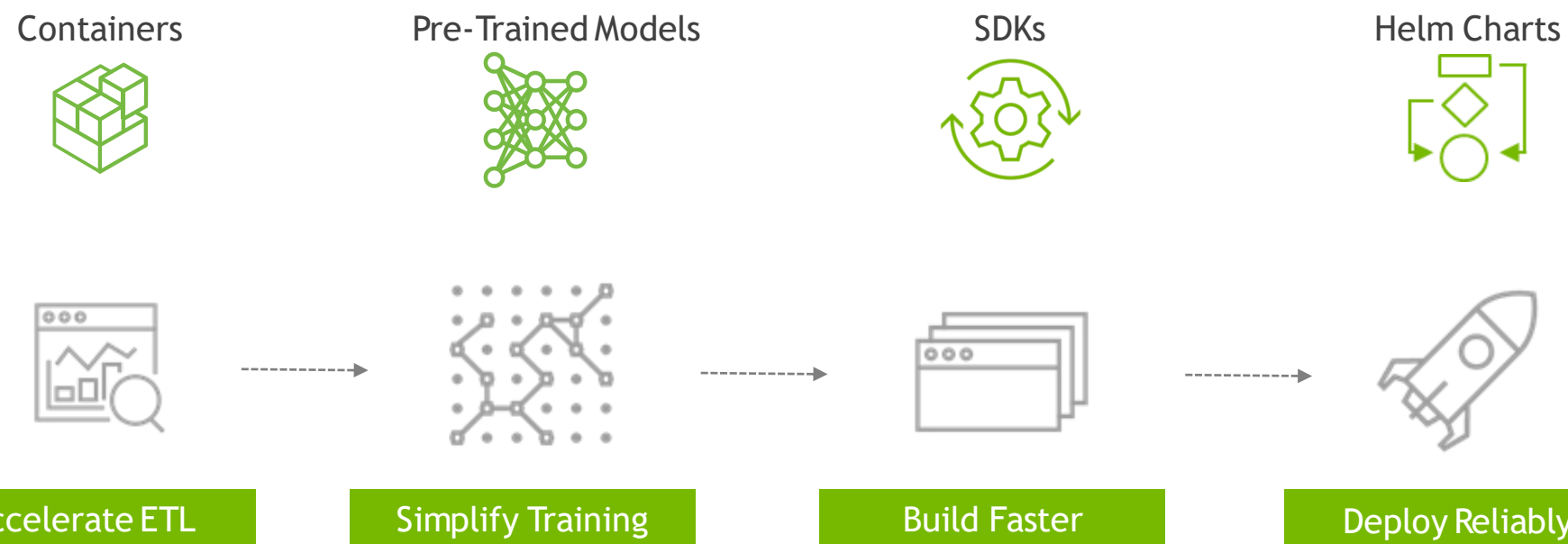# AI WORKFLOWS ARE COMPLEX

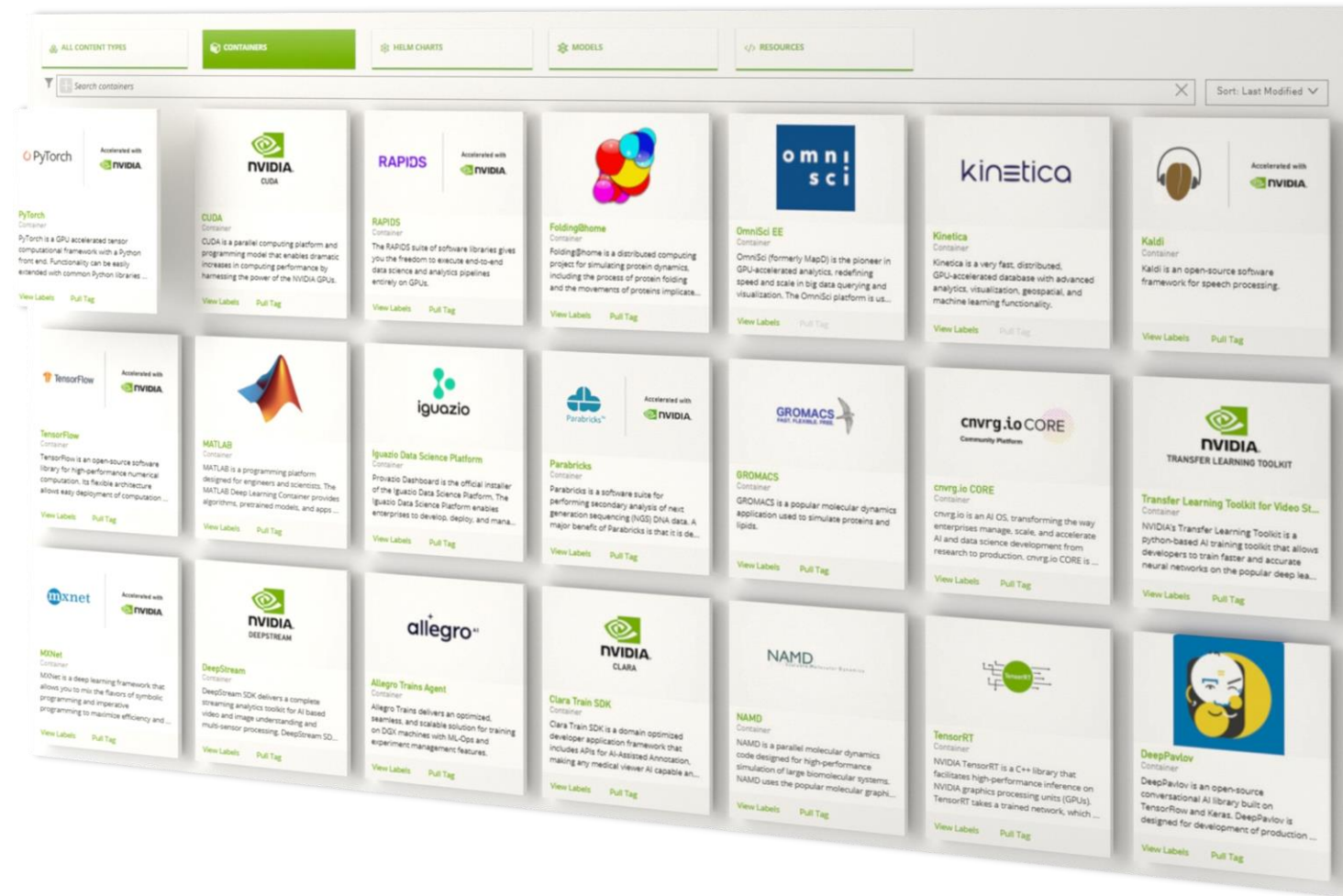| Data Engineer | | | | | Data Scientist | | | | | Developer | | | | DevOps |

Extract · Transform · Load · Feature Engineering → Setup · Build Model · Train · Optimize → Integrate · Optimize · Build App → Deploy & Monitor

Timeline

# NGC CATALOG HELPS SIMPLIFY AND ACCELERATE AI WORKFLOWS

| Data Engineer | Data Scientist | Developer | DevOps |
|---|---|---|---|

Extract    Transform    Load    Feature Engineering    Setup    Build Model    Train    Optimize    Integrate    Optimize    Build App    Deploy & Monitor

Timeline

Containers    Pre-Trained Models    SDKs    Helm Charts

Accelerate ETL    Simplify Training    Build Faster    Deploy Reliably

# NGC CONTAINERS ENABLE YOU TO FOCUS ON BUILDING AI



## ENTERPRISE READY SOFTWARE

Scanned for CVEs, malware, crypto

Tested for reliability

Backed by Enterprise support

## PERFORMANCE OPTIMIZED

Scalable

Updated Monthly

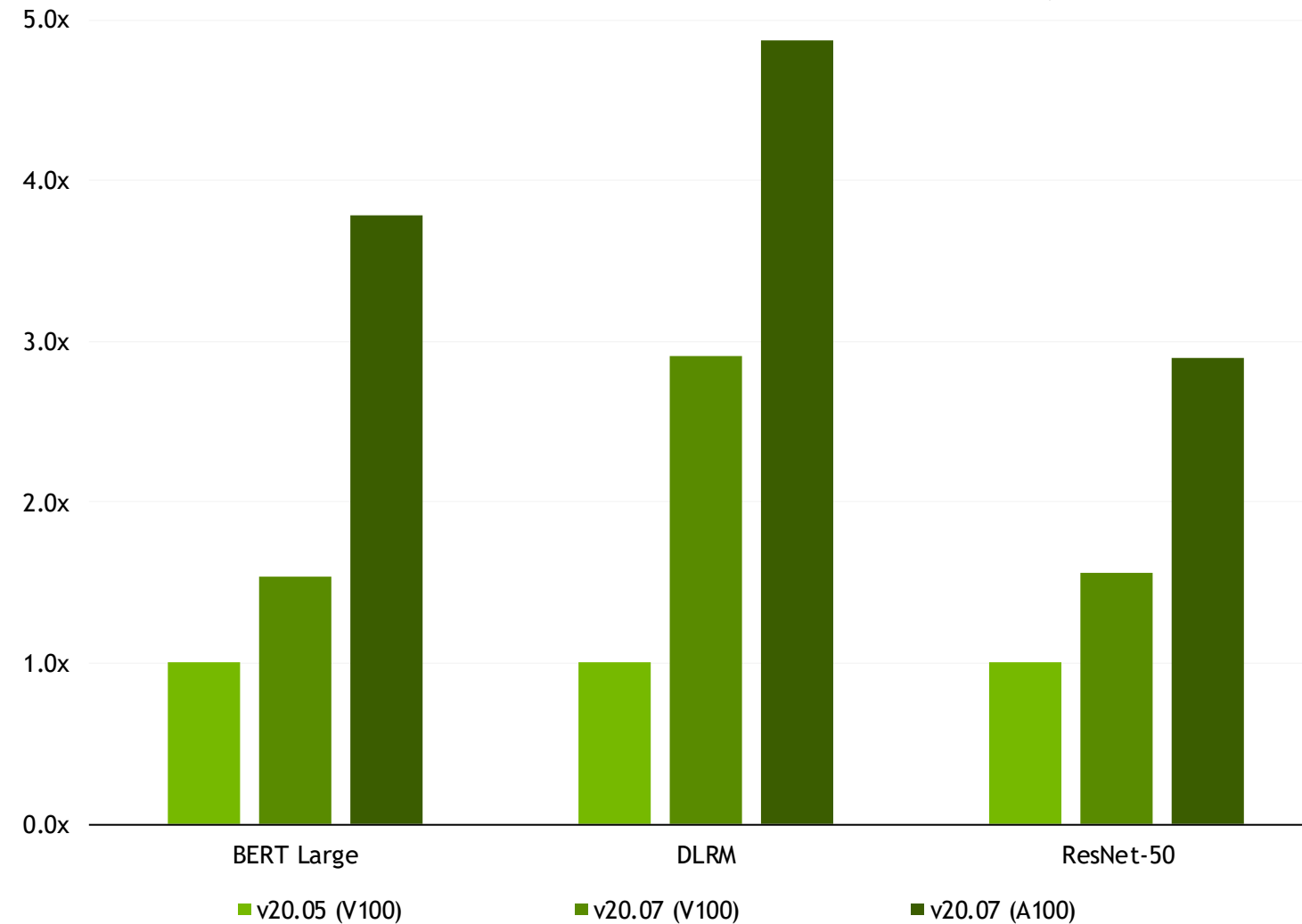Better performance on the same system

## DEPLOY ANYWHERE

Docker, cri-o runtimes

Bare metal, VMs, Kubernetes

Multi-cloud, on-prem, hybrid, edge

# DO WHAT YOU DO BEST, FASTER



**ENTERPRISE READY SOFTWARE**

Scanned for CVEs, malware, crypto

Tested for reliability

Backed by Enterprise support

**PERFORMANCE OPTIMIZED**

Scalable

Updated Monthly

Better performance on the same system

**DEPLOY ANYWHERE**

Docker, Singularity runtimes

Bare metal, VMs, Kubernetes

Multi-cloud, on-prem, hybrid, edge

BERT-Large and ResNet-50 v1.5 Training performance with TensorFlow on a single node 8x V100 (32GB) & A100 (40GB). Mixed Precision. Batch size for BERT: 10 (V100), 24 (A100), ResNet: 512 (V100, v20.05), 256 (v20.07)
DLRM Training performance with PyTorch on 1x V100 & 1x A100. Mixed Precision. Batch size 32768. DRLM trained with v20.03 and v20.07

# EASILY IDENTIFY THE RIGHT MODELS WITH CREDENTIALS



## WIDE RANGE OF USE CASES

ResNet-50, SSD, MobileNet, VGG16

WaveGlow, BERT, NeMo

Wide & Deep, DLRM & many more

## PRE-TRAINED MODELS

Faster training

Higher accuracy
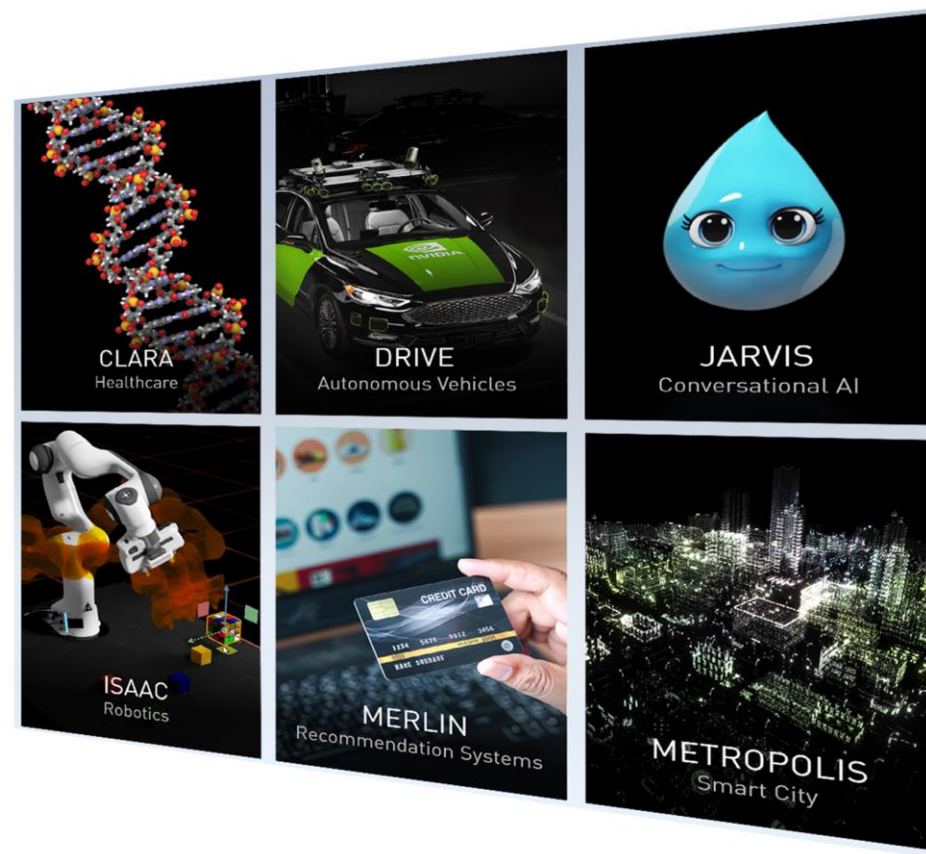
Transparency through credentials

## RESOURCES

Get started with code samples

Customize NGC models

Reproduce with recipes

# INDUSTRY APP FRAMEWORKS FOR END-TO-END AI WORKFLOWS



**TRANSFER LEARNING TOOLKIT**

Domain adaptability

Significantly reduce development time

**TENSORRT**

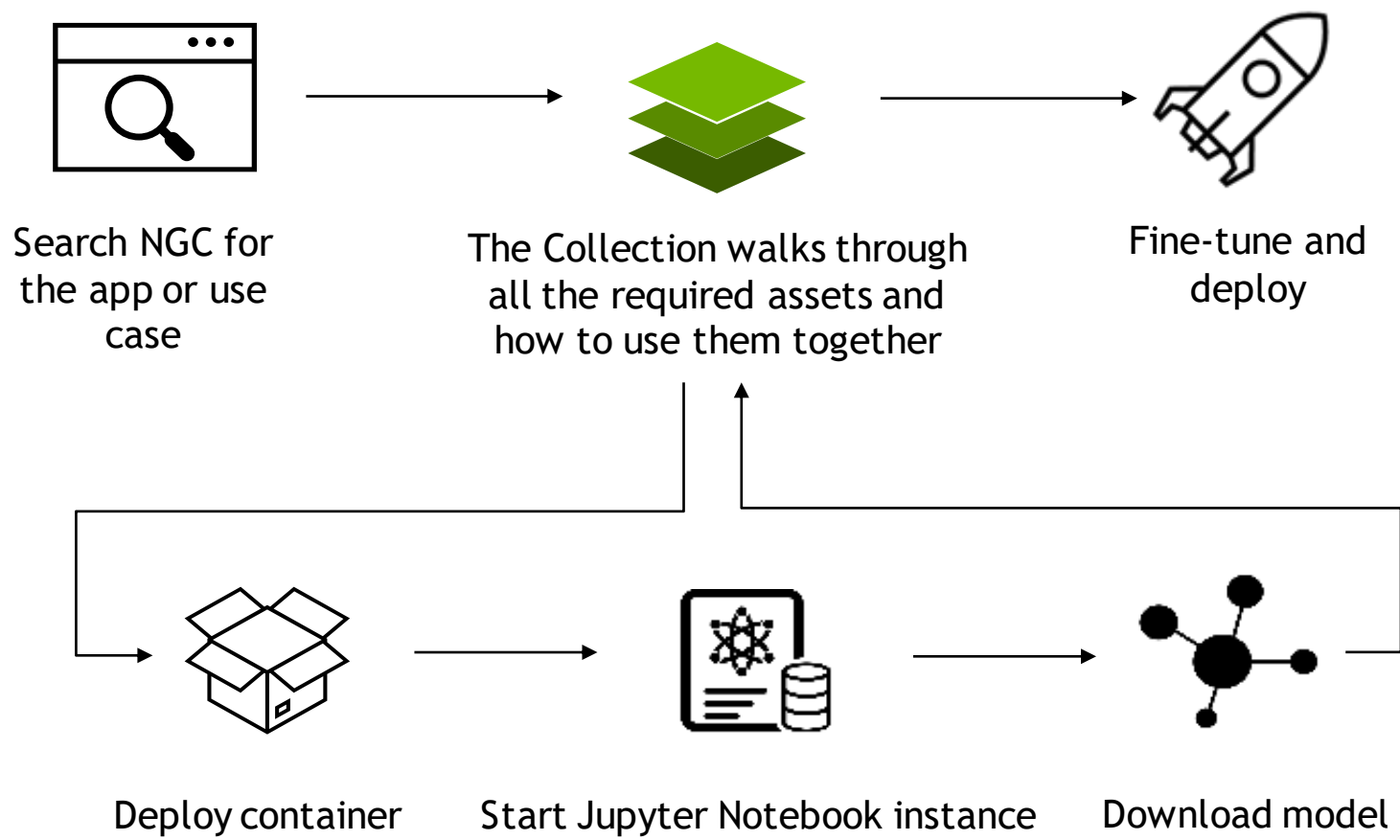Optimizes for low latency and high-throughput

Integrated with major frameworks

**TRITON**

High performance inference on GPU/CPU systems

Supports multiple frameworks backends

# EVERYTHING YOU NEED TO BUILD YOUR AI IN ONE LOCATION



Search NGC for the app or use case

The Collection walks through all the required assets and how to use them together

Fine-tune and deploy

Deploy container

Start Jupyter Notebook instance

Download model

## COLLECTIONS

Compatible assets grouped together, removes guesswork

Curated software by use cases

Detailed documentation further simplifies work for users

## READY-TO-USE

Conversational AI

Computer Vision

NVIDIA AI App Frameworks
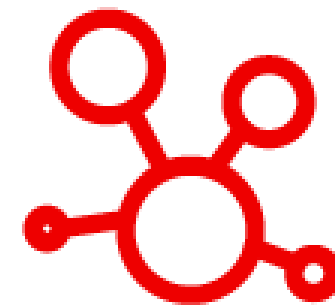
# WHY KUBERNETES AND DEVOPS FOR AI/ML?

**AGILITY**

Respond quickly
with automated compute resource
management, and increased
collaboration

**CONSISTENCY & PORTABILITY**

Develop and deploy ML models
consistently across data center,
edge, and public clouds.

**FLEXIBILITY**

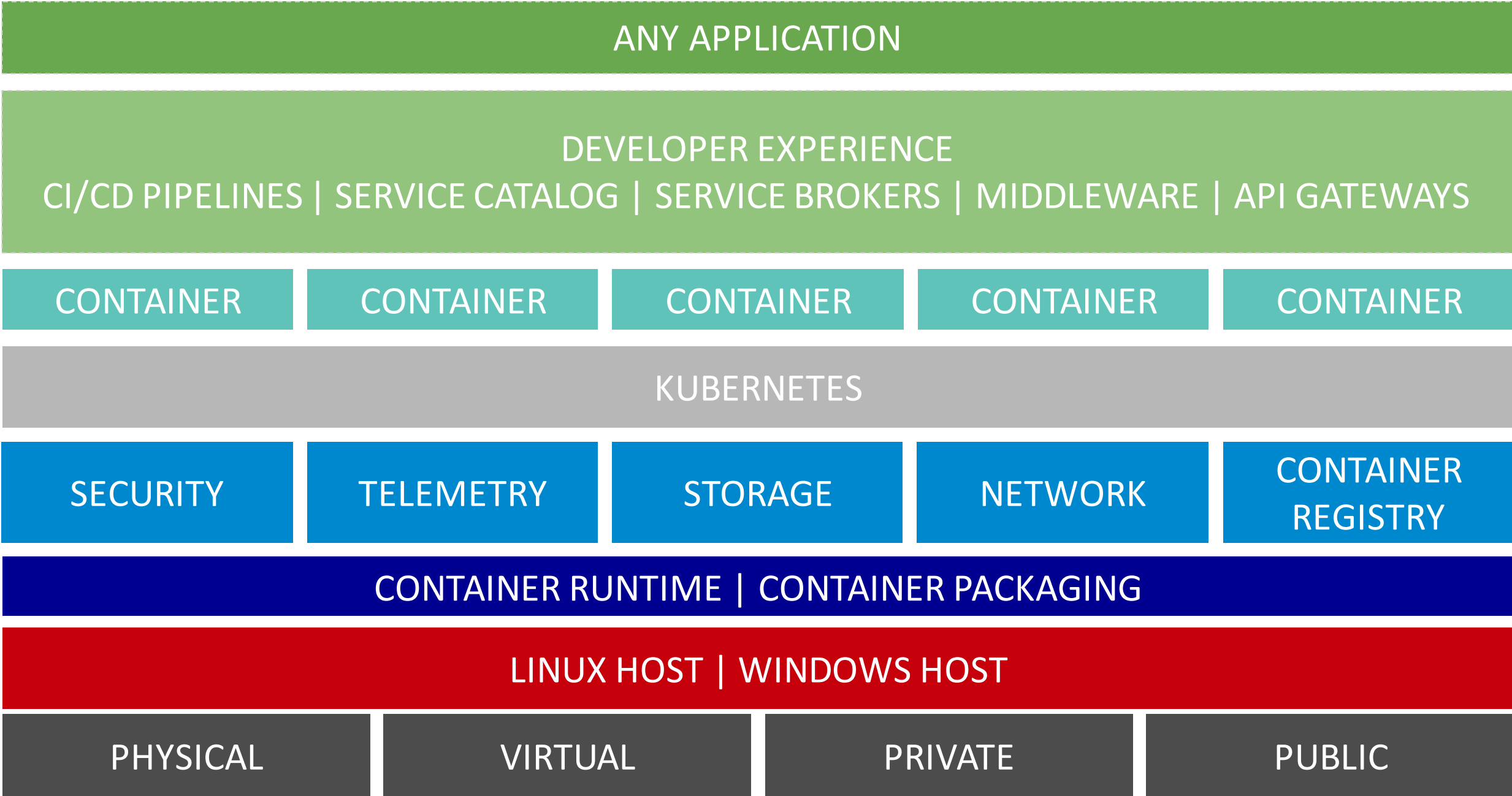Provision AI/ML environments as and
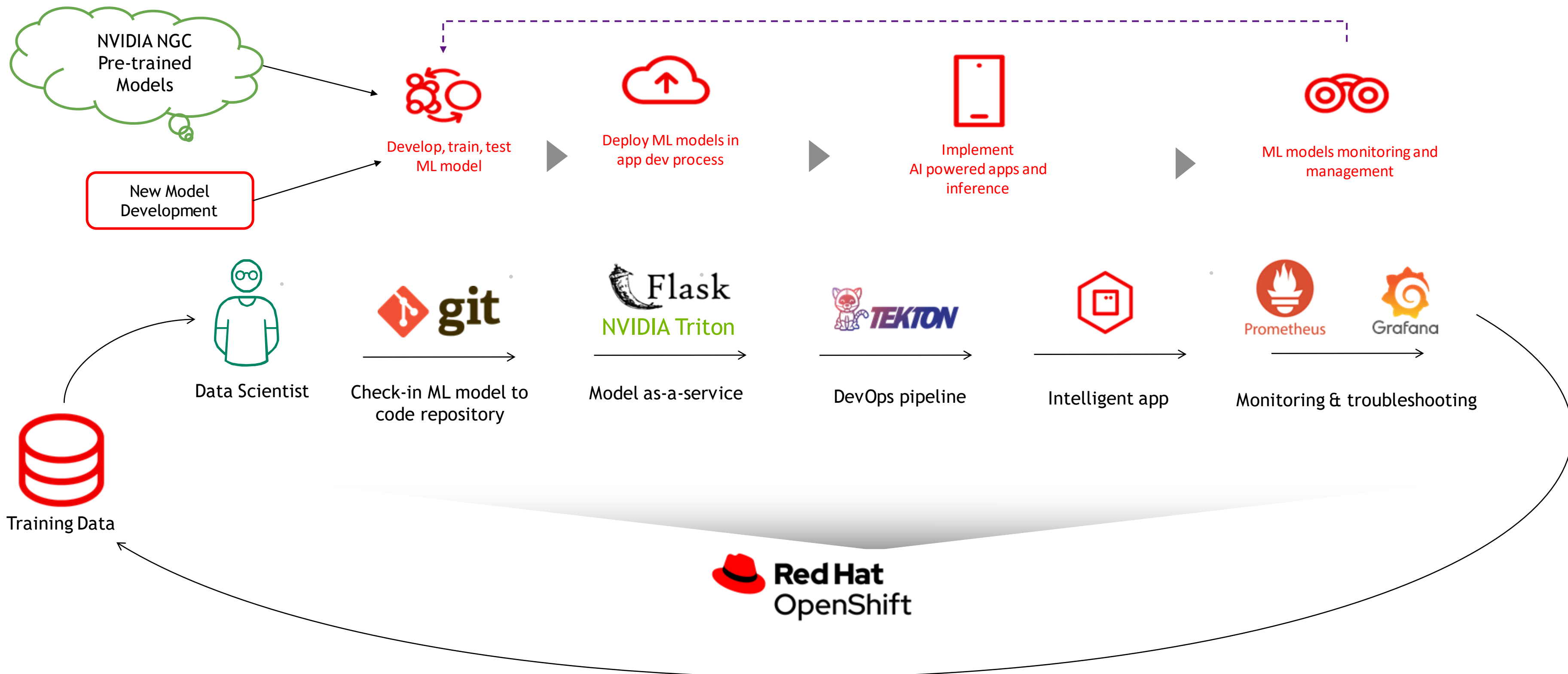when you need them.

**SCABILITY**

Autoscaling and high availability of
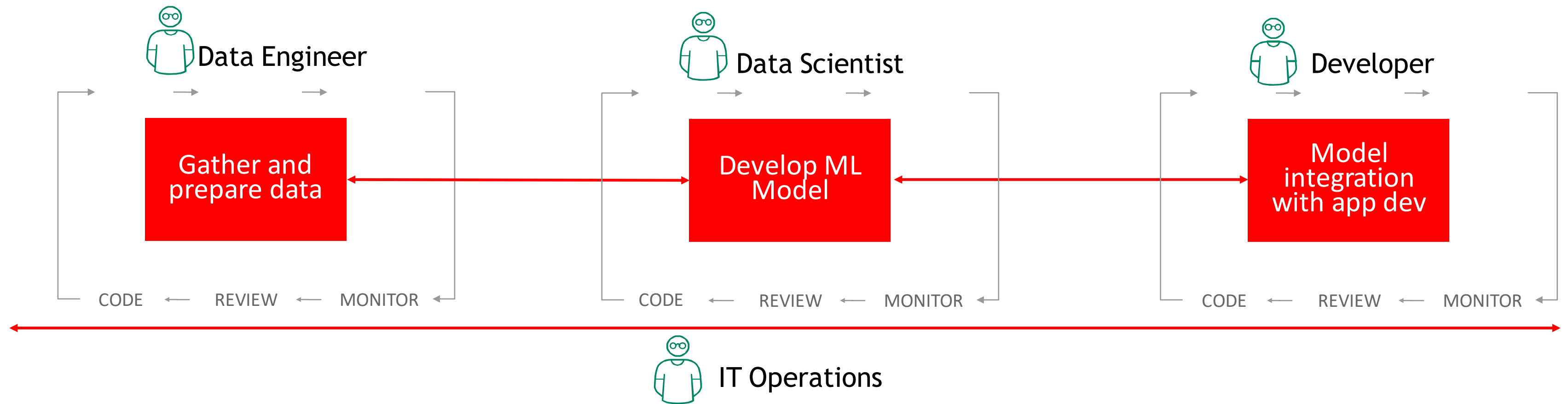the AI/ML solution stack.

# OPENSHIFT ENABLES CONTAINERS, KUBERNETES, AND DEVOPS IN PRODUCTION

**Red Hat OpenShift**

| ANY APPLICATION | | | | |
|---|---|---|---|---|
| DEVELOPER EXPERIENCE<br>CI/CD PIPELINES \| SERVICE CATALOG \| SERVICE BROKERS \| MIDDLEWARE \| API GATEWAYS | | | | |
| CONTAINER | CONTAINER | CONTAINER | CONTAINER | CONTAINER |
| KUBERNETES | | | | |
| SECURITY | TELEMETRY | STORAGE | NETWORK | CONTAINER REGISTRY |
| CONTAINER RUNTIME \| CONTAINER PACKAGING | | | | |
| LINUX HOST \| WINDOWS HOST | | | | |
| PHYSICAL | VIRTUAL | PRIVATE | PUBLIC | |

**NVIDIA.**   **Red Hat**

# RED HAT OPENSHIFT HELPS FAST TRACK AI/ML LIFECYCLE

NVIDIA NGC Pre-trained Models

New Model Development

Develop, train, test ML model

Deploy ML models in app dev process

Implement AI powered apps and inference

ML models monitoring and management

Data Scientist

**git**
Check-in ML model to code repository

**Flask**
**NVIDIA Triton**
Model as-a-service

**TEKTON**
DevOps pipeline

Intelligent app

**Prometheus** **Grafana**
Monitoring & troubleshooting

Training Data

**Red Hat OpenShift**

**NVIDIA.** **Red Hat**

# RED HAT OPENSHIFT HELPS FAST TRACK AI/ML LIFECYCLE

Data Engineer

Data Scientist

Developer

Gather and prepare data

Develop ML Model

Model integration with app dev

CODE    REVIEW    MONITOR

CODE    REVIEW    MONITOR

CODE    REVIEW    MONITOR

IT Operations

**Red Hat OpenShift**

**Container, Kubernetes and DevOps Platform**

- File (NFS, HDFS), Object (S3) and Block
- High Throughput, Low Latency, Secure
- Data Movement - Kafka
- Data Analytics  - Spark and BDC
- Data pipelines  - Tekton and ArgoCD

- CPUs, memory, GPUs, FPGA
- High speed networking and storage
- Containers with - Language (python), frameworks (PyTorch),  IDE (Jupyter)
- On-demand scale up
- DevOps - Build v2, Tekton, ArgoCD

- RESTful services for models
- Services monitoring and alerting
- Services logging and diagnostics

NVIDIA.   Red Hat

CONVERSATIONAL AI DEMO

# TASK: TRAIN, OPTIMIZE & DEPLOY A FINE-TUNED BERT MODEL

## Demo Workflow

Training

Deployment



1. Load training data on OpenShift Container Storage

2. Get BERT container from NGC

3. Customize training yaml and mount persistent volume

4. Train on OpenShift using A100 or V100 or T4 GPUs

5. Optimize model using TensorRT

6. Get Triton container from NGC

7. Load on OpenShift Container Storage
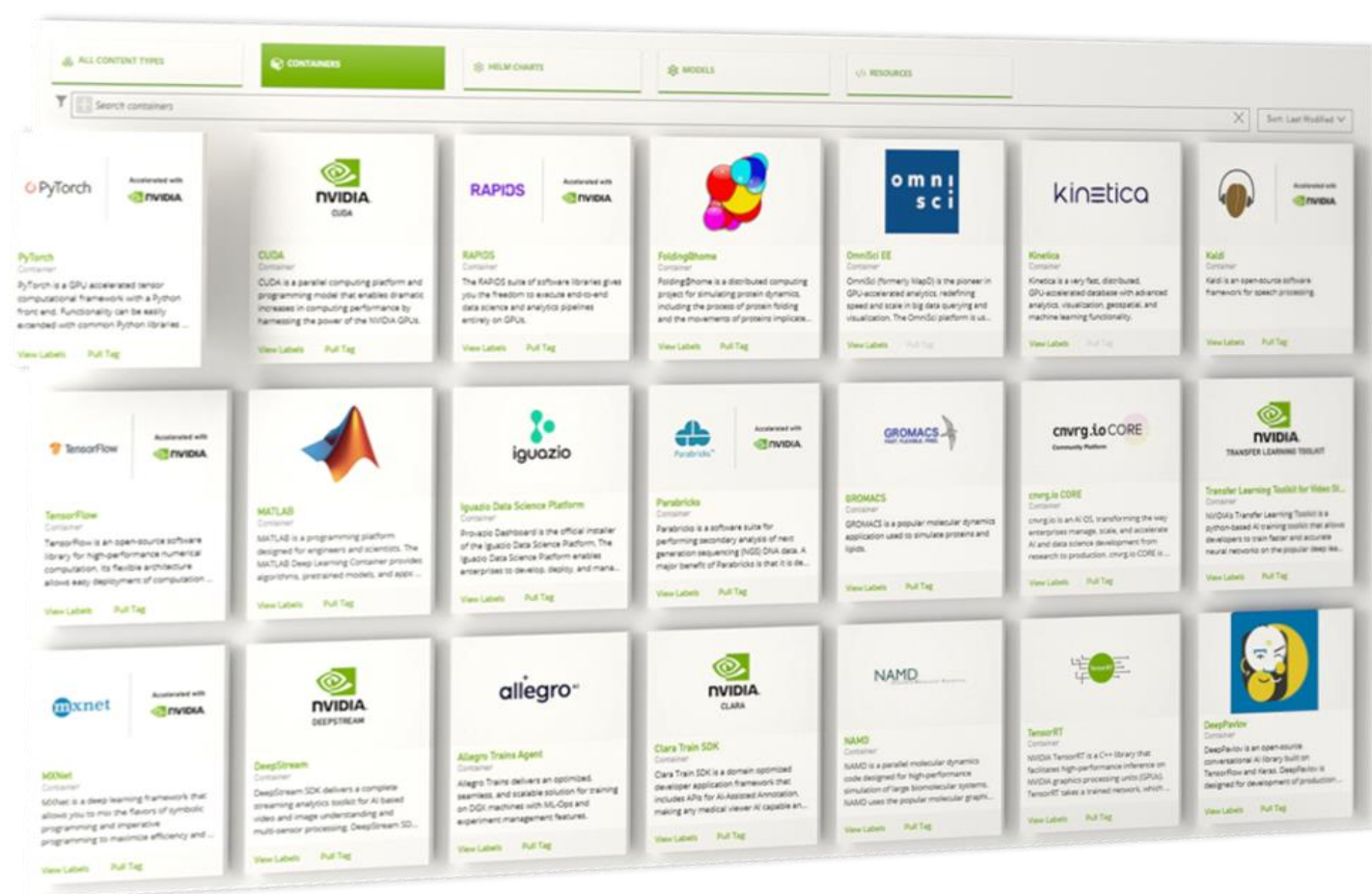
8. Deploy on OpenShift using A100/V100 /T4 GPUs

LET'S SEE IT IN ACTION

# BUILD AI FASTER WITH NGC AND OPENSHIFT

## Pull the OpenShift Collection from NGC and Run on OpenShift

### ngc.nvidia.com

### openshift.com/try | openshift.com/nvidia