# Case Study: Machine Learning as a Service in Production

November 20th, 2019
San Diego

**Itay Gabbay
Machine Learning
Lead, MOD (Israel)**

**Tushar Katarki
Product Manager,
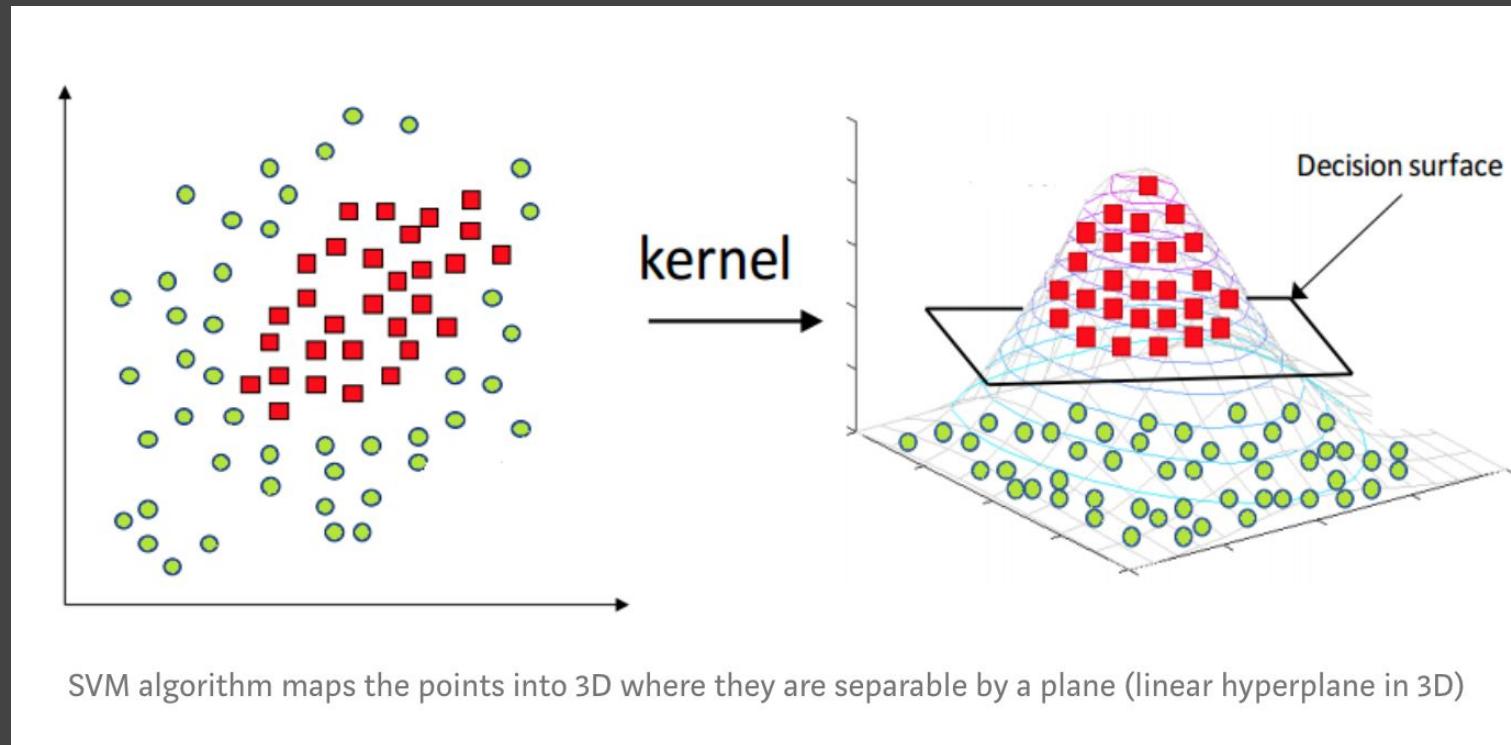ML on OpenShift at
Red Hat**

Digital Transformation

Private Cloud

Accelerate R&D Projects

Machine Learning

# Why Machine Learning?

# Machine Learning vs. Traditional Methods



SVM algorithm maps the points into 3D where they are separable by a plane (linear hyperplane in 3D)
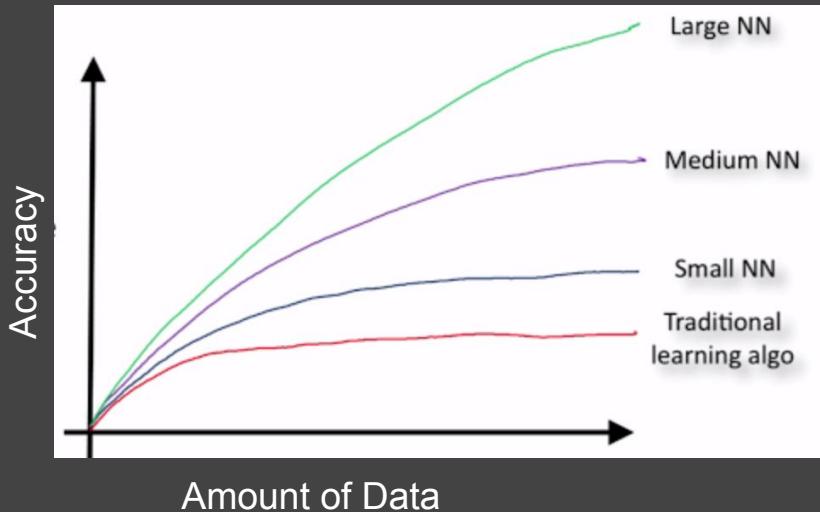
# Deep Learning
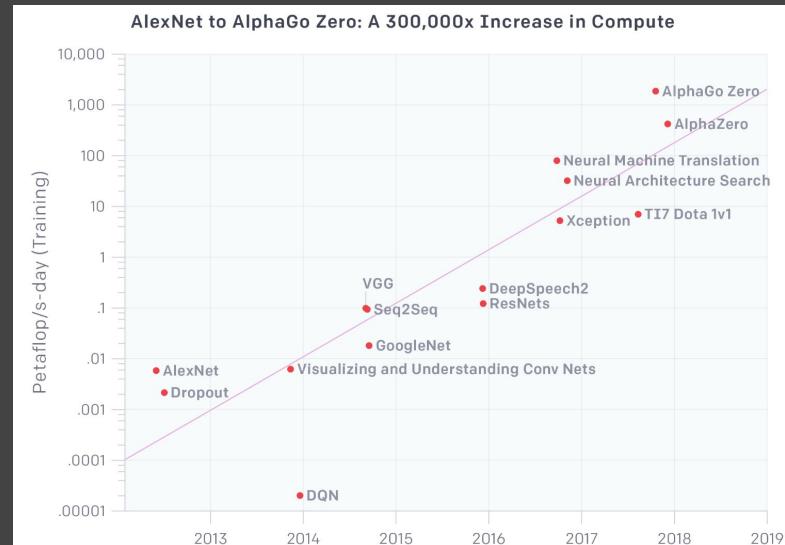
# AI vs Humans !



Go champion Lee Sedol, on the right, concedes the second of possible five games vs. Google's AlphaGo AI.

# Characteristics of Machine Learning
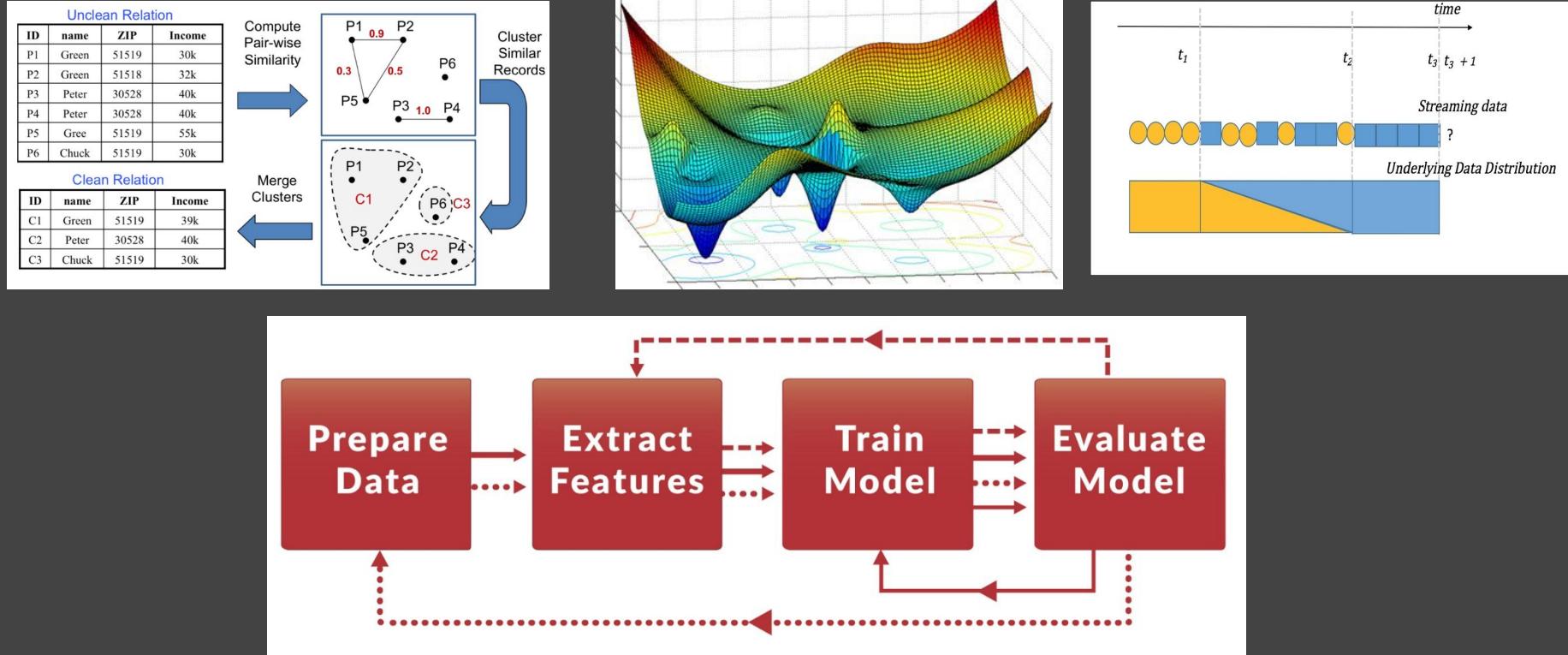
# Data and Compute Intensive



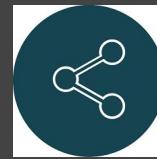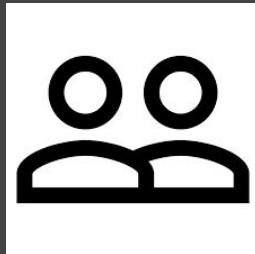Source: Andrew Yan-Tak Ng, Chief Scientist at Baidu Research



Source: OpenAI

# ML is Iterative !

# Collaboration and Sharing

# The Journey @MOD
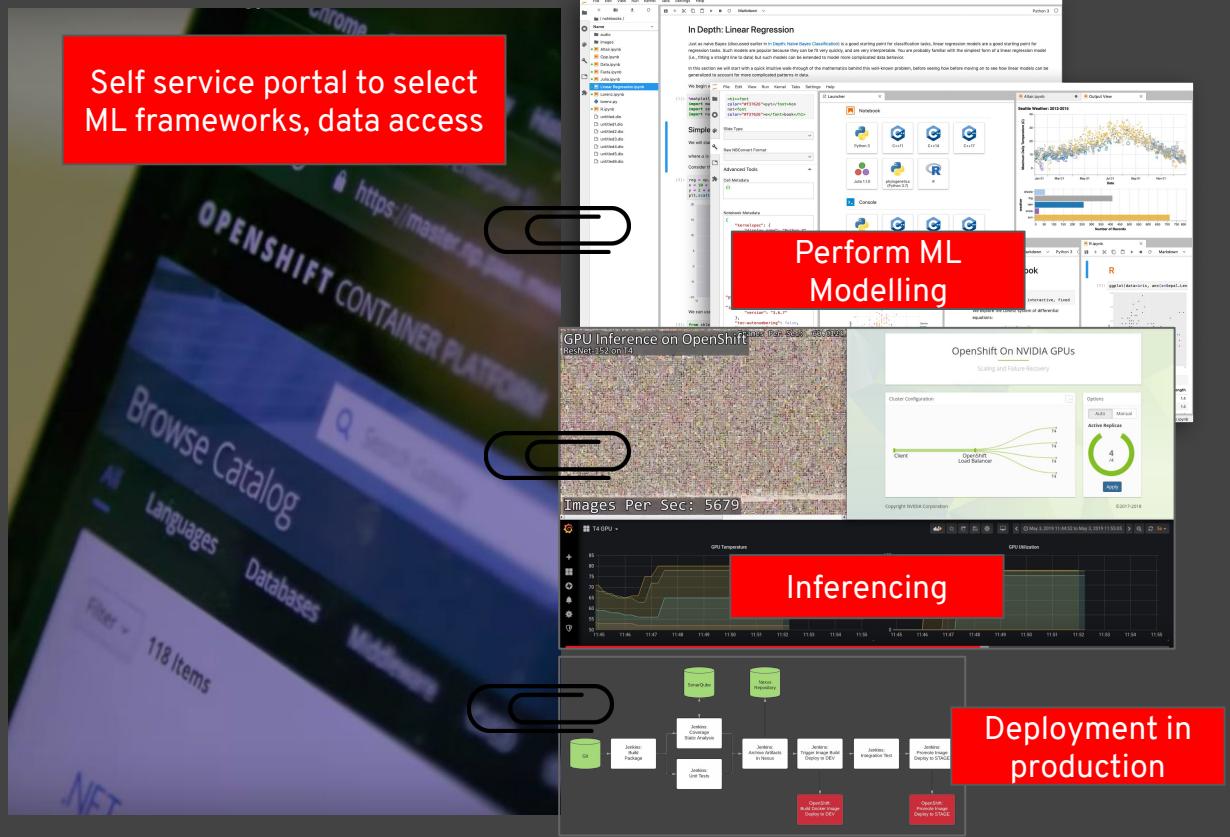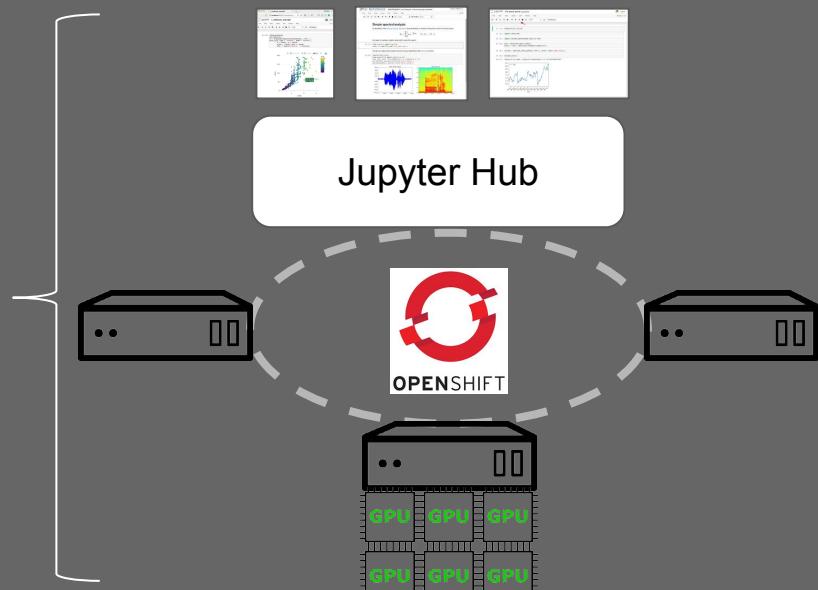
# Top 3 Requirements



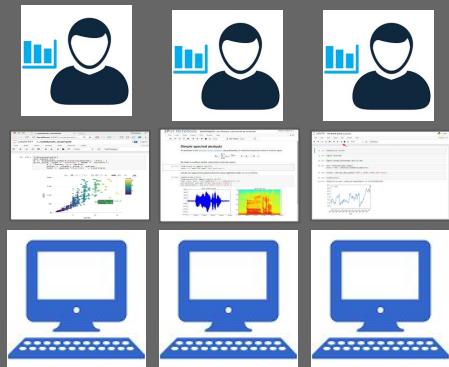Access to elastic compute



Deploy models into production

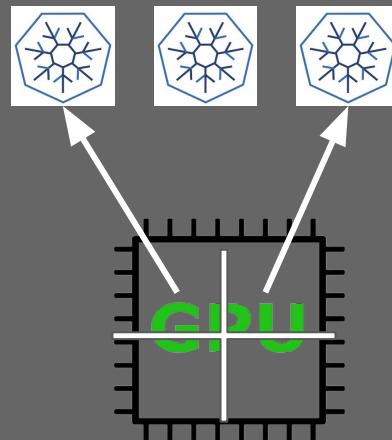

Optimize and Scale Machine Learning

As a Data Scientist, I want a *"self-service cloud like"* experience for my Machine Learning projects, where I can access a **rich set of modelling frameworks, data, and computational resources,** *share and collaborate* with colleagues, and deliver my work into **production** with **speed,** *agility and repeatability to* **drive organizational value**!

Self service portal to select ML frameworks, data access

Perform ML Modelling

Inferencing

Deployment in production

# Step 1: A Self Service Cloud
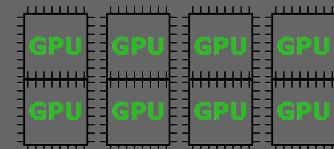
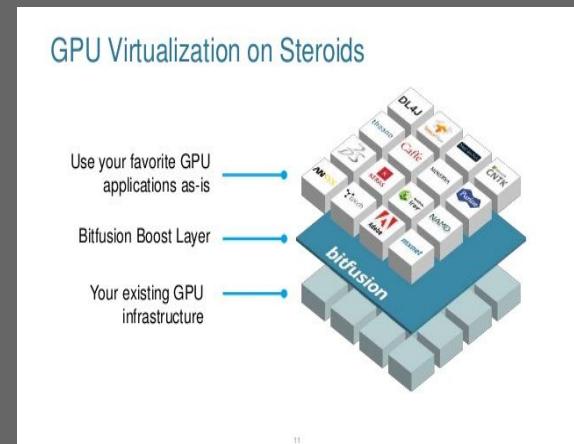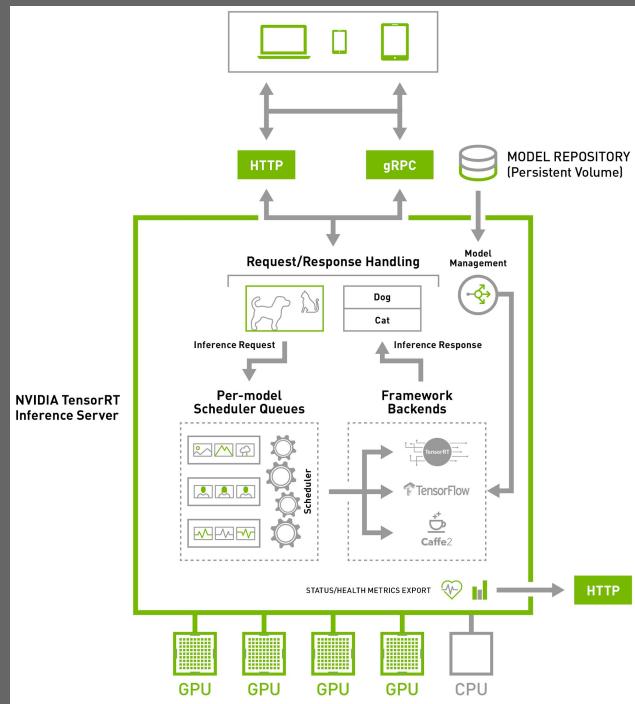# Multiple GPUs and Sharing GPUs

# FYI - Other GPU Sharing techniques

# Multi Tenancy



- → Identity and Role Based Access
- → Network Policies
- → Storage Classes
- → Resource Isolation
- → Resource Quotas

# Step 2: Scaling ML Talent

WHAT IF I TOLD YOU

YOU CAN BE A ML EXPERT

# Introducing AutoML

# AutoML Design Elements

# 30%

Average improvement in machine learning model performance

# 70%

Increase in machine learning models number

# 600%

Growth of users in the previous 6 months

# Challenges and Next Steps

➔ **Automate** development, debug and deployment of **notebooks**

➔ Better way to **save and catalog** experiments

➔ AutoML for unstructured data - **images, audio**

➔ Supported way for **GPU sharing**

➔ Multi-cluster

# What's Next

# Community First



NVIDIA NGC



ML-as-a-service reference architecture on OpenShift and open source and ISV content

Home for k8s community to share operators for various apps/tools

# Welcome to OperatorHub.io

OperatorHub.io is a new home for the Kubernetes community to share Operators. Find an existing Operator or list your own today.

CATEGORIES

85 ITEMS                                          VIEW ▦ ⌄    SORT A-Z ⌄

AI/Machine Learning
Application Runtime
Big Data
Cloud Provider
Database
Developer Tools
Integration & Delivery
Logging & Tracing
Monitoring
Networking
OpenShift Optional
Security
Storage
Streaming & Messaging

PROVIDER

☐ Altinity (1)
☐ Amazon Web Services (1)
☐ Anchore (1)
☐ Appsody (1)
☐ Aqua Security (1)
Show 57 more

CAPABILITY LEVEL

☐ Basic Install (40)
☐ Seamless Upgrades (13)
☐ Full Lifecycle (18)
☐ Deep Insights (12)
☐ Auto Pilot (2)

### Akka Cluster Operator
provided by Lightbend, Inc.

Run Akka Cluster

### Altinity ClickHouse Operator
provided by Altinity

ClickHouse Operator

### Anchore Engine Operator
provided by Anchore Inc.

Apache CouchDBâ„¢ is a highly available

### Apache CouchDB
provided by IBM

Apache CouchDBâ„¢ is a highly available

### Apache Spark Operator
provided by radanalytics.io

An operator for
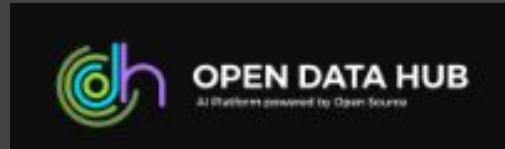
### Appsody Operator
provided by Appsody

Deploys Appsody based applications

### Aqua Security Operator
provided by Aqua Security, Inc.

The Aqua Security

### AtlasMap Operator
provided by AtlasMap

AtlasMap is a data mapping solution with

### AWS S3 Operator
provided by Red Hat

Manage the full lifecycle of installing, configur

### AWS Service Operator
provided by Amazon Web Services, Inc.

The AWS Service

### Banzai Cloud Kafka Operator
provided by Banzai Cloud

### Camel K Operator
provided by The Apache Software Foundation

Apache Camel K is a

### Cassandra
provided by Instaclustr

Manage the full lifecycle of the Cassandra

### CockroachDB
provided by Helm Community

CockroachDB Opera

### Community Jaeger Operator
provided by CNCF

Provides tracing,

### Crunchy PostgreSQL Enterprise
provided by CrunchyData.com

Install full-stack

### Dynatrace OneAgent
provided by Dynatrace LLC

Install full-stack

### Eclipse Che
provided by Eclipse Foundation

A Kube-native

### Elastic Cloud on Kubernetes
provided by Elastic

Run Elasticsearch,

### EnMasse
provided by EnMasse

EnMasse provides messaging as a

# Open Data Hub

# OpenShift 4

**Operator**
based installer

→

Build, Event and Serve with Knative and Tekton

OpenShift Service Mesh
(Istio + Jaeger + Prometheus + Kiali)

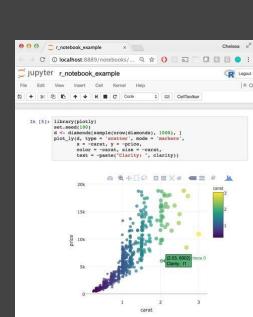OpenShift Container Platform
(Enterprise Kubernetes)

GPU

Datacenter

Cloud

# From experimentation to production with CI/CD
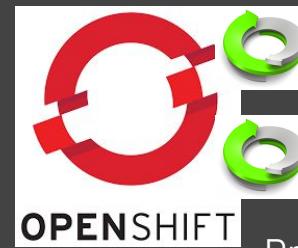


Check-in to source repo

Source-2-image

Deloy notebook container

Model test & iteration and integration

Container

OPENSHIFT

Promote and Serve models into production as services

Continuous monitoring for performance and drift

# DEMO - Self Service with Open Data Hub

# SUMMARY

➔ MOD Case Study: Machine Learning-as-a-service platform

   ◆ Why and how they built a **cloud-like** experience

   ◆ AutoML

➔ Kubernetes and OpenShift and open source tools

➔ OperatorHub and OpenDataHub

# THANK YOU !

➔ Contact:

➔ [tkatarki@redhat.com](mailto:tkatarki@redhat.com)

➔ @tkatarki

➔ Q/A