

# Using Kubernetes to Host Machine Learning Workloads in a Cloud Native Way

And what you need to know to get there.....

Michael Doherty

Technical Architect, WW DC Team

DGTL-BRKCLD-2148



#CiscoLive



# Topics Covered

- Opportunity knocks!
- So how do WE view Machine Learning?
- How and where is Machine Learning done?
- So where does kubernetes & kubeflow fit in?
- Cisco ML Solutions
- Opportunities to look out for
- The new frontier....The Edge!

The background is a dark blue field filled with numerous small, semi-transparent squares and dots in various colors including light blue, teal, yellow, orange, and red. These elements are scattered across the frame, with a higher concentration of yellow and orange squares forming a diagonal streak from the top right towards the bottom right.

Opportunity Knocks!

# So let's set a baseline

Humans learn from experience



Computers follow instructions



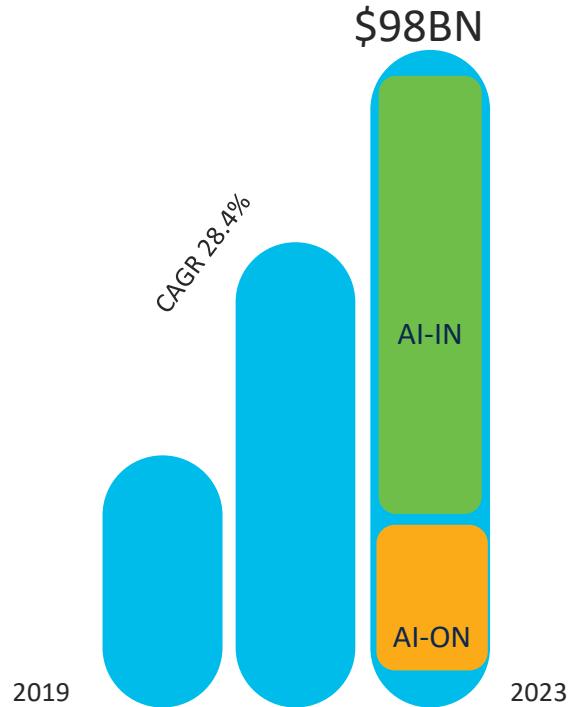
AI/ML learns from data





- AI set to be the biggest technology transition in history.
- We do have a role to play here.
- AI sales enablement will yield greater returns in AI.

# The AI Opportunity - Cisco



Source: IDC AI systems spending guide 2019

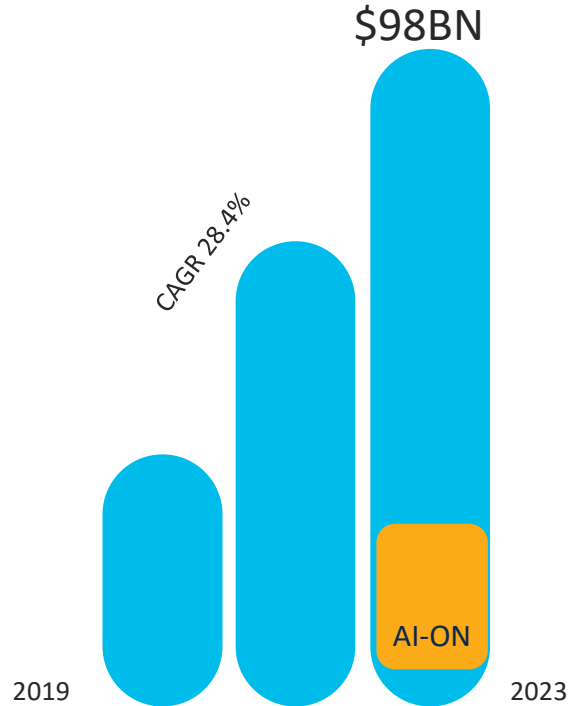


Infrastructure for AI  
Platforms  
DC | IoT | Meraki



AI Powered Infrastructure  
Features  
DNAC | Collab | Security

# The AI Opportunity - Customers



Source: IDC AI systems spending guide 2019



25%

of a  
**Data Scientist** time  
is spent on  
infrastructure tasks.

[MIT Technology Review](#)



“..majority of teams  
developing ML  
capabilities are doing so  
using **open-source**  
tooling because of the  
**dearth of viable  
commercial options**”.

[Gartner](#)

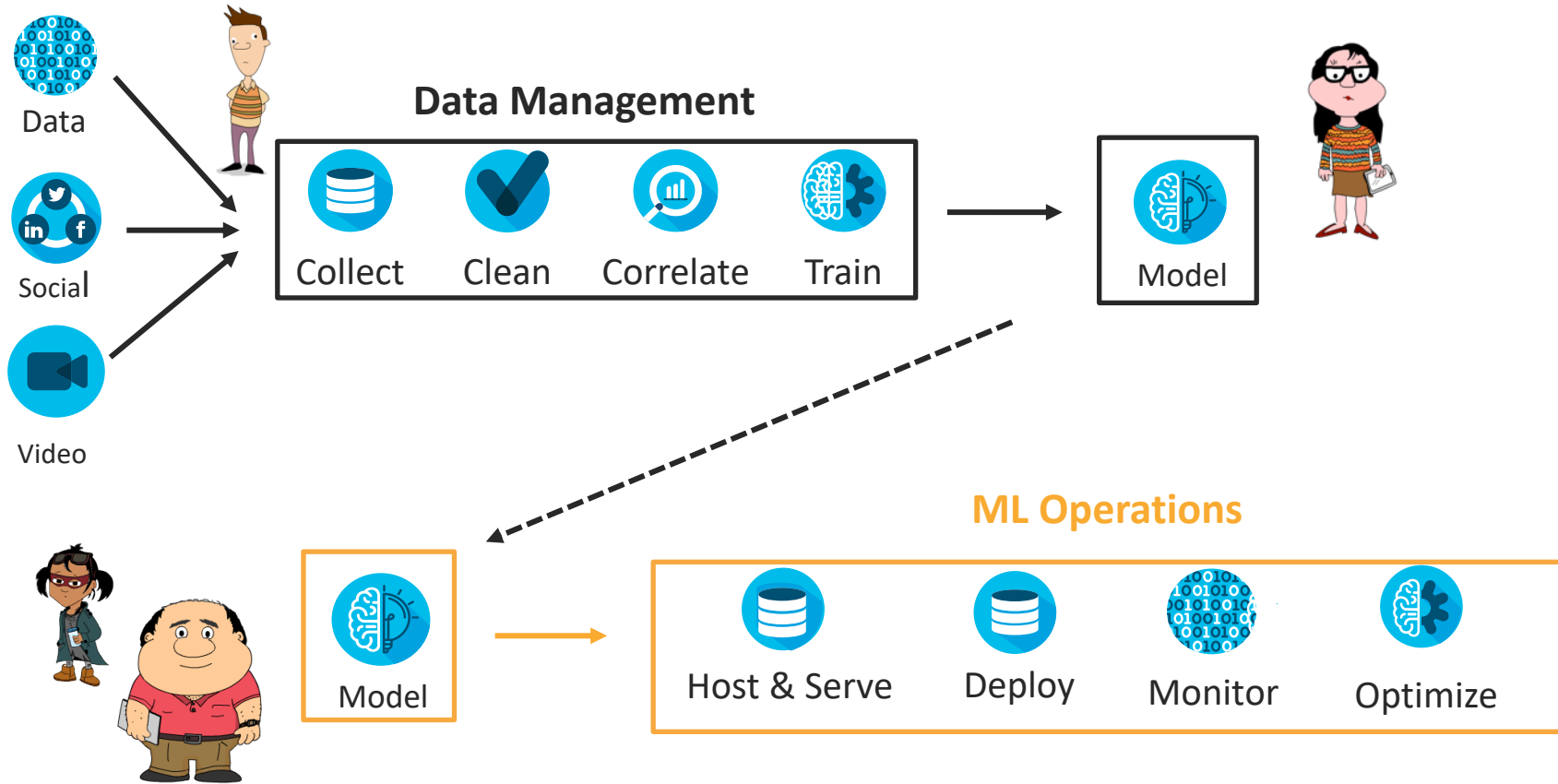


The background is a dark blue field filled with numerous small, semi-transparent squares and dots in various colors including light blue, teal, yellow, orange, and red. These elements are scattered across the frame, with a higher concentration of yellow and orange squares forming a diagonal streak from the top right towards the bottom right.

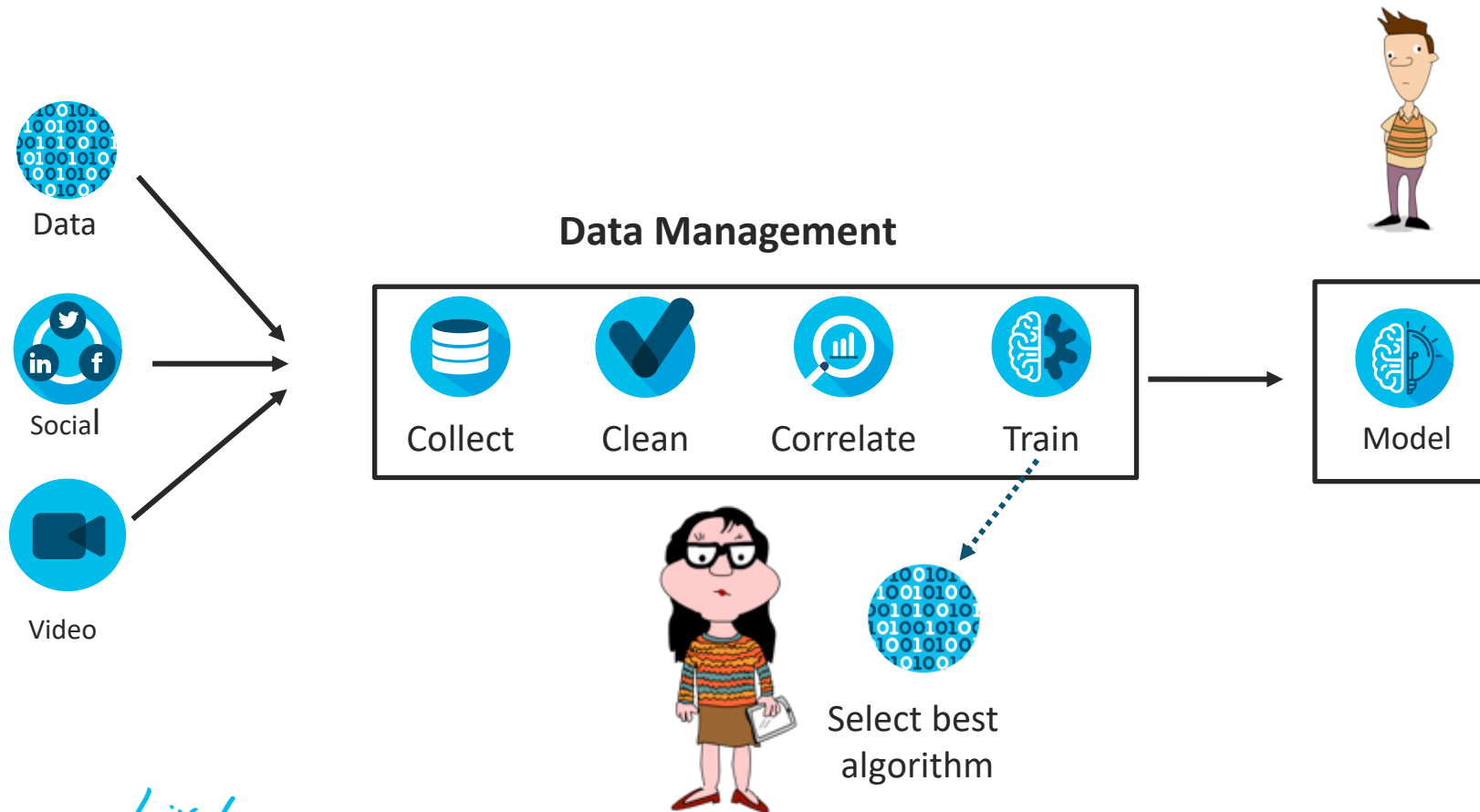
So how do WE see  
Machine Learning?



# There 2 main parts to consider

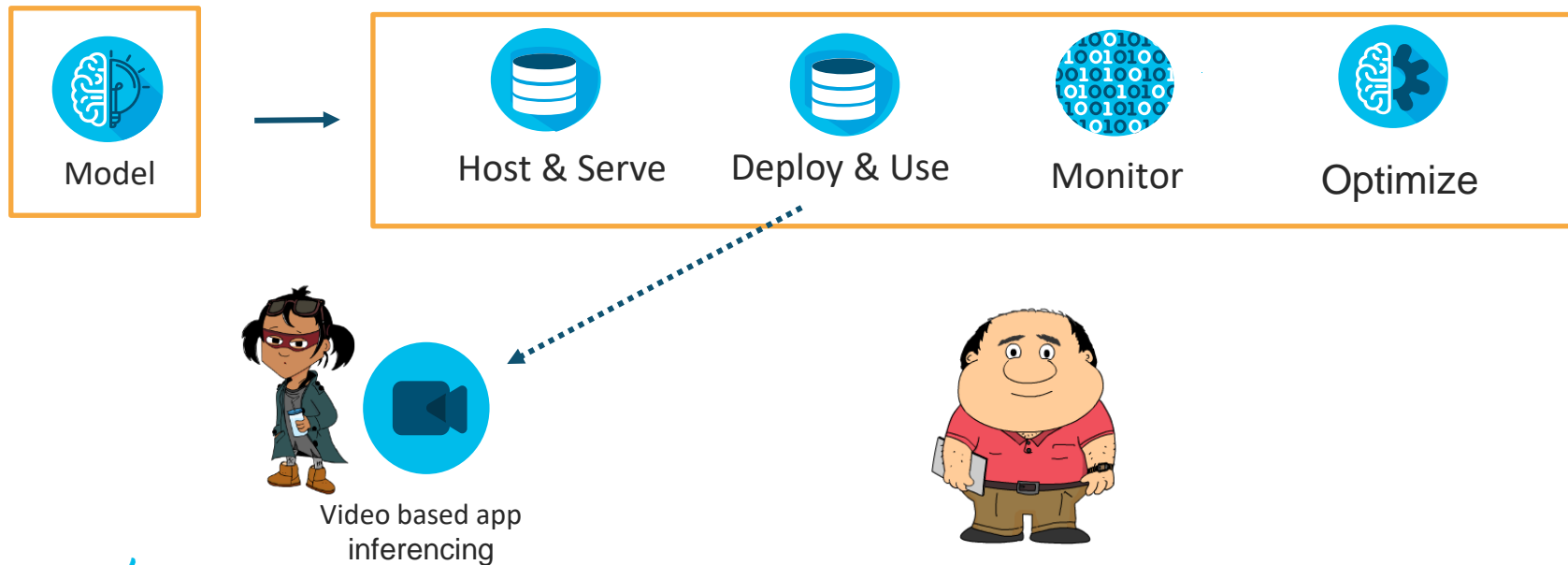



# The Data Scientists/Engs



# The Operations Teams

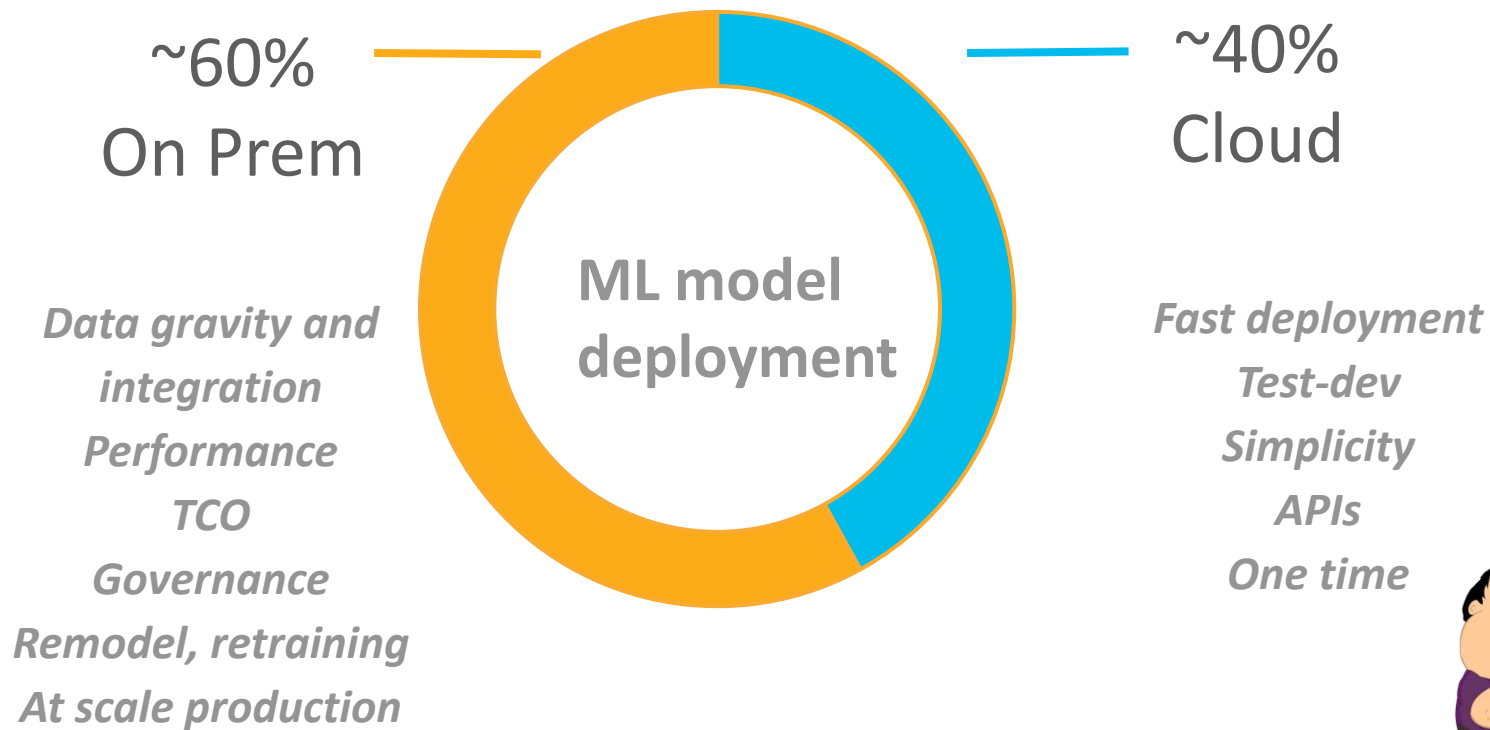
## ML Operations



The background is a dark blue field filled with numerous small, semi-transparent squares and dots in various colors including light blue, teal, yellow, orange, and red. These elements are scattered across the frame, with a higher concentration of yellow and orange squares forming a diagonal streak from the top right towards the bottom right.

# How and where is Machine Learning done?

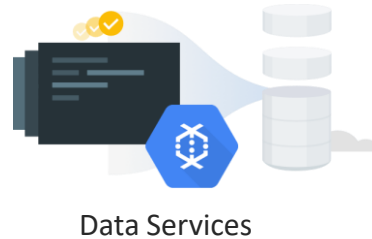
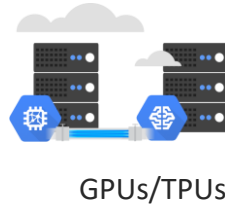
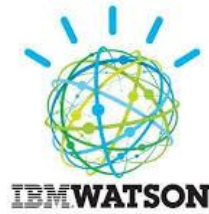
# On-Premise vs. Cloud



Source: Gartner, Market Guide for Machine Learning Compute Infrastructures, Sep 2108, ID: G00362287; Figure 3



# Public Cloud Appeal



**CISCO** *Live!*



# Private Cloud Drivers



## Costs @ Scale



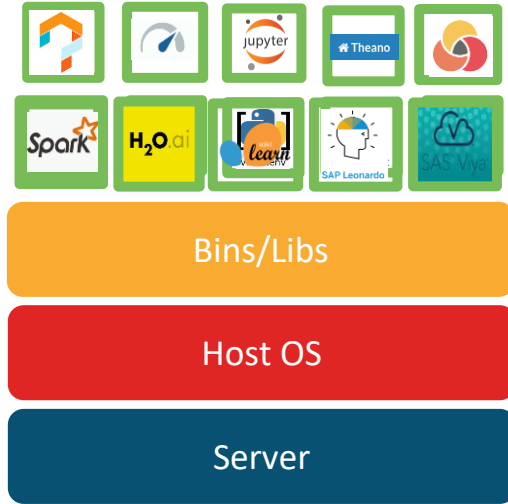
## Data Privacy



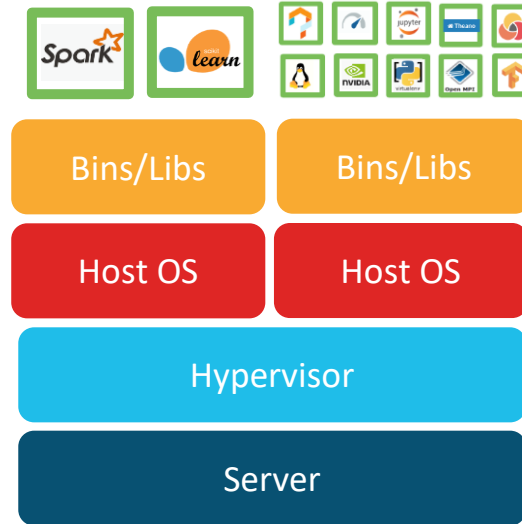
## Data gravity and integration



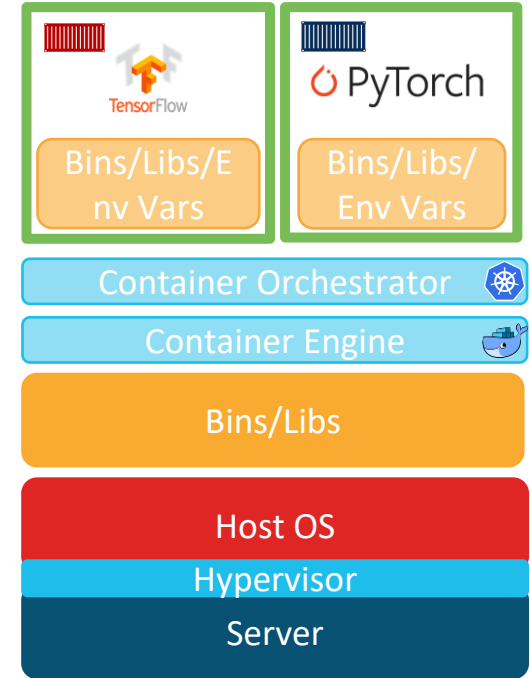
# Traditional vs Cloud Native



Bare Metal



Virtual Machines

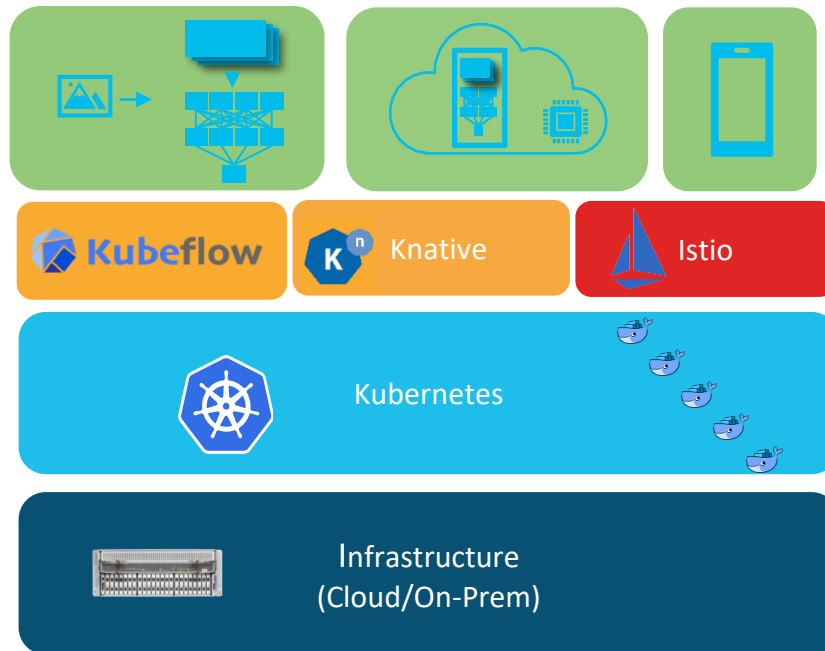


Containers

The background is a dark blue field filled with numerous small, semi-transparent squares and dots in various colors including light blue, teal, yellow, orange, and red. These elements are scattered across the frame, with a higher concentration of yellow and orange squares forming a diagonal streak from the top right towards the bottom right.

So where does  
kubernetes & kubeflow  
fit in?

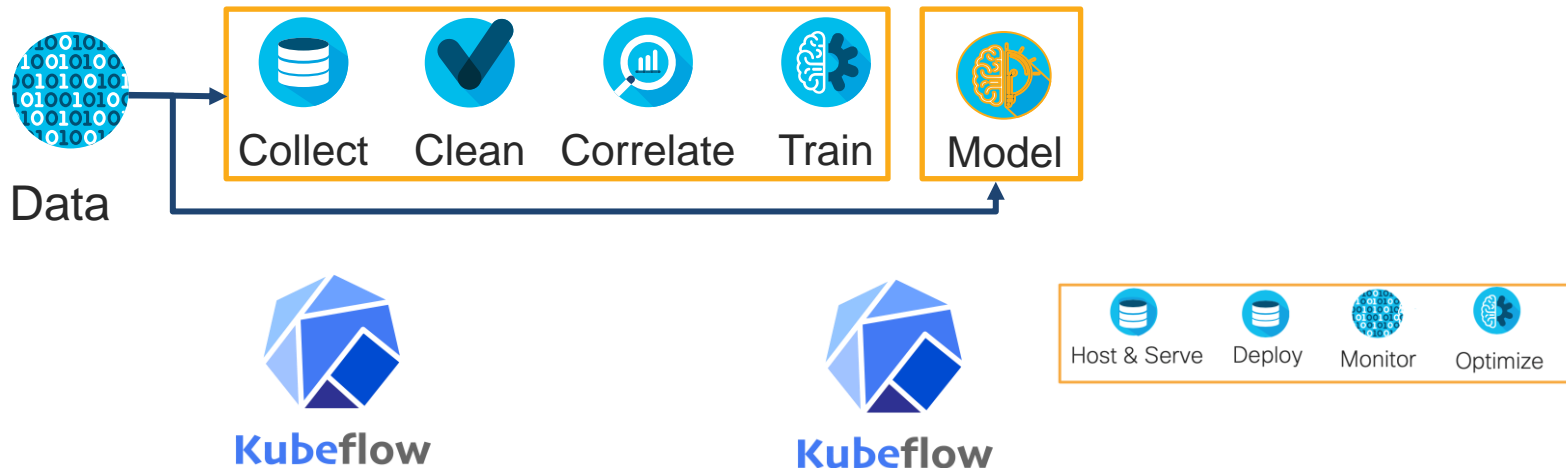
# K is the magic letter...well almost!



# Kubeflow



*Machine learning lifecycle manager that makes it easy to develop, deploy and manage portable, scalable end-to-end ML workflows everywhere*



# Kubeflow

## Scalable ML Services on Kubernetes



### Easy to get started

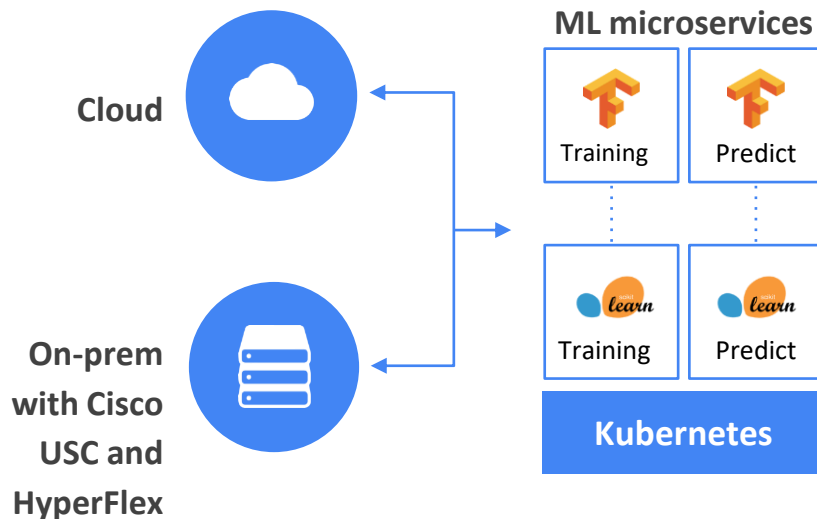
- Out-of-box support for top frameworks
  - pytorch, caffe, tf and xgboost
- Kubernetes manages dependencies, resources

### Swappable & Scalable

- Library of ML Services
- GPU support
- Massive Scale

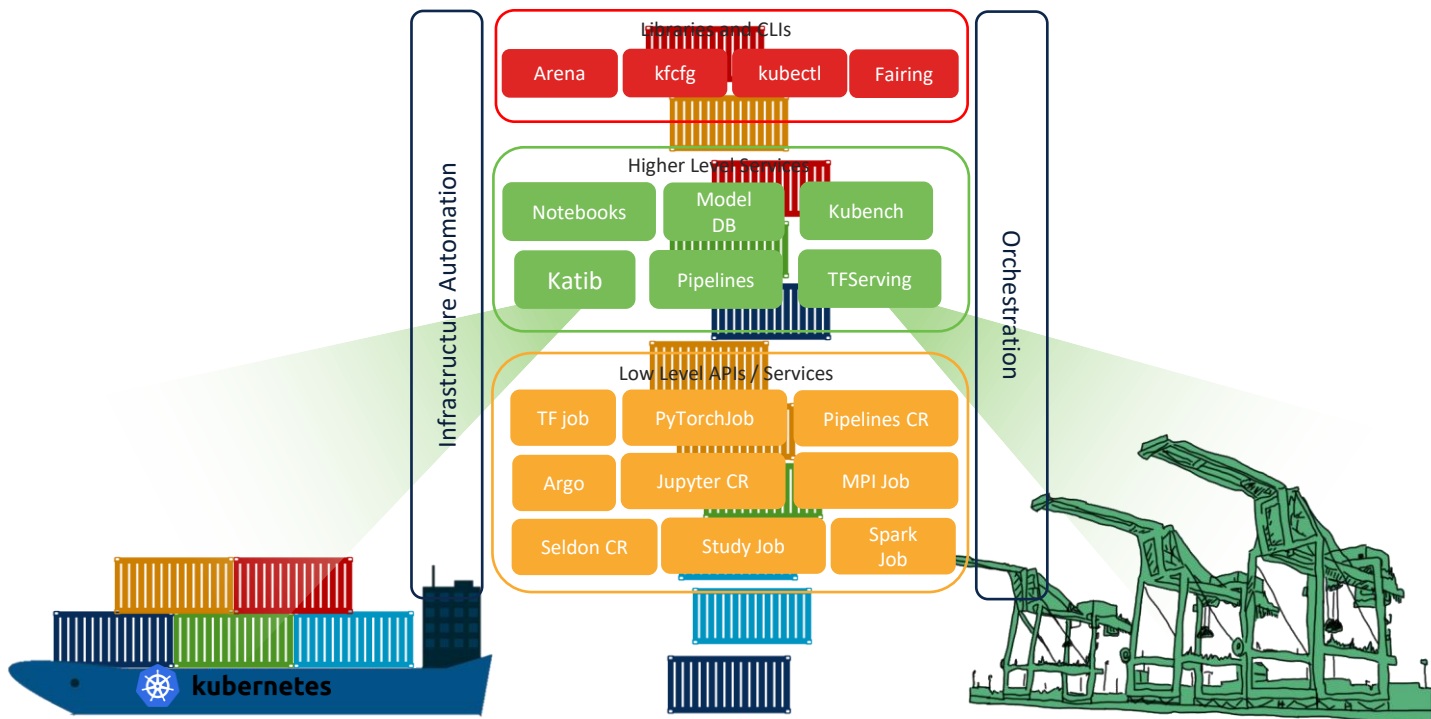
### Meet customer where they are

- Public Cloud
- On-prem





# Kubeflow is made up of multiple components



# High Level Design PoV for Kubeflow Pipeline

## Training a Neural Net to recognize logos

### Inputs (collect)

~2k images of logos scraped from the web, labelled "Cisco" or "Not-Cisco"



### Pre-processing (clean)

- add noise & other perturbations
- rescale to 400x300px
- hold out some data for validation.



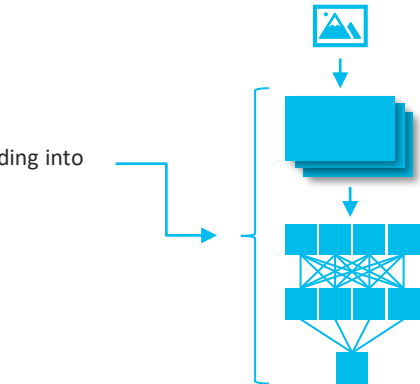
### Model Specification

- Convolutional "funnel" architecture feeding into densely connected binary classifier.
- Specification of hyperparameters

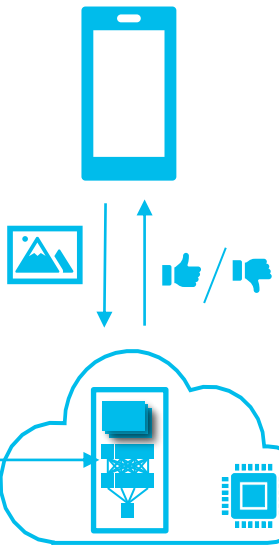
### Model Training

### Model Evaluation

### Export trained model



## Deploying the N.N. in an app



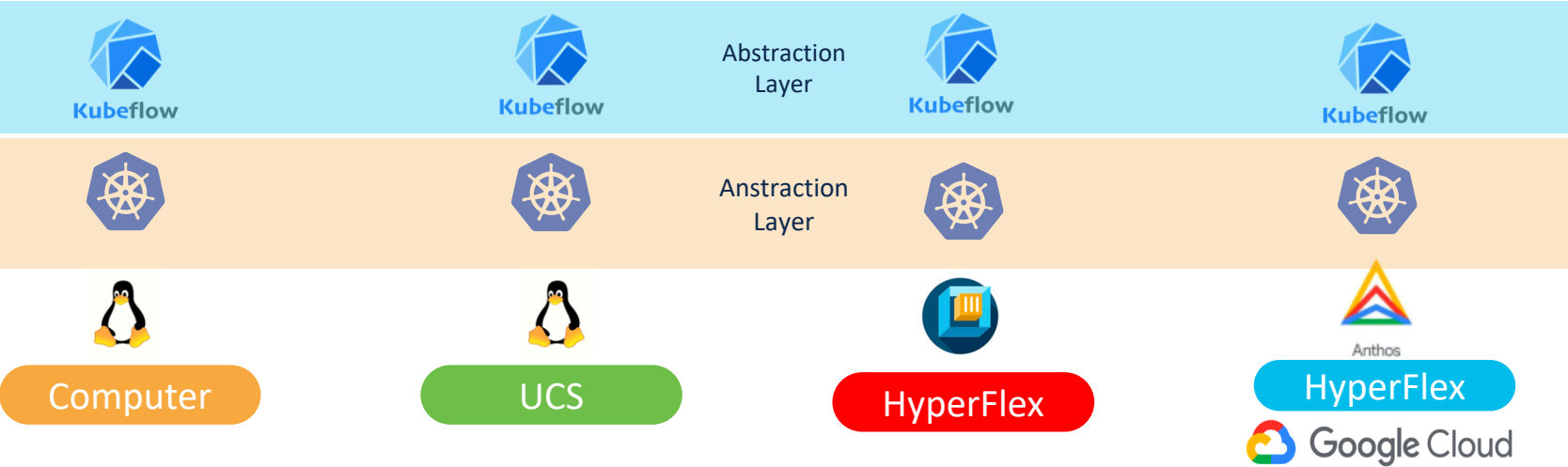
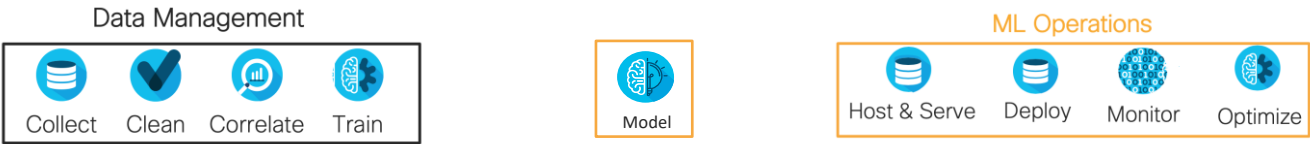
Frontend

API

Backend

Repeat

# Consistent Cloud and On-Premise Machine Learning Experience

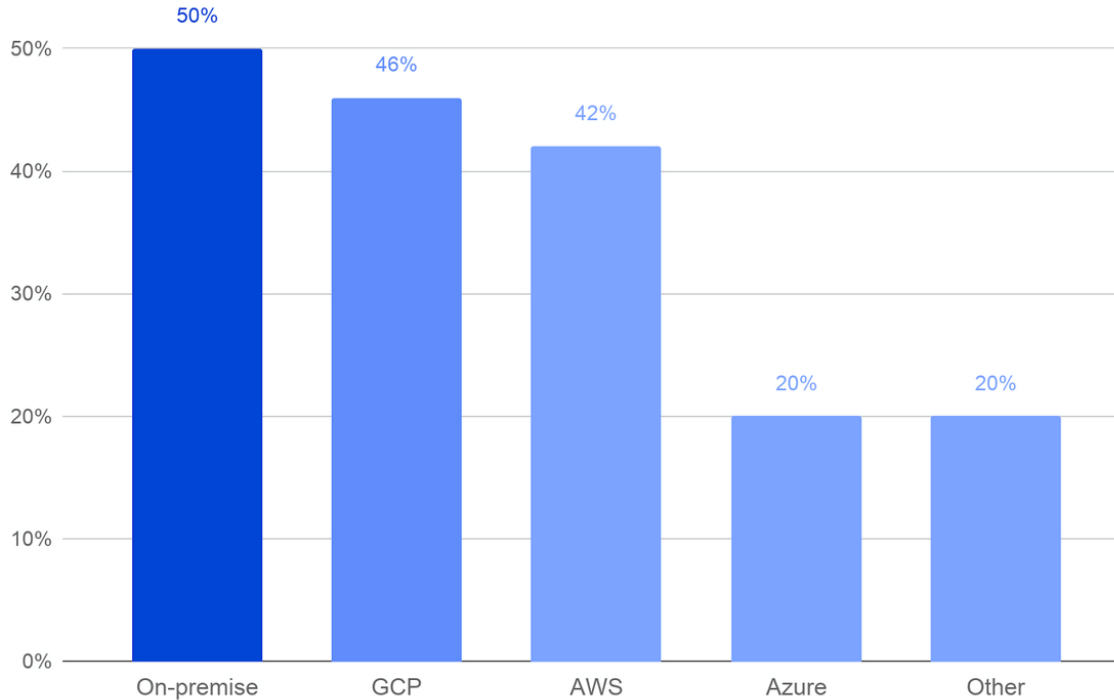


# Demo Intro



*“Data Scientists have to wait on average 3 months to get an ML Platform delivered to them by IT”*

# Combined with.....



Pic. 3. Where do you run your AI/ML workloads? (Multiple select)? N = 50



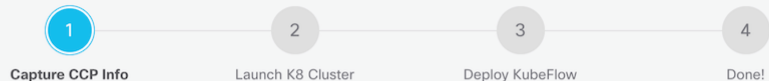
# Introducing ML Anywhere

A ML platform deployment 'easy button' built on Kubernetes and Kubeflow



ML Anywhere Deployment

Logging



CCP IP Address

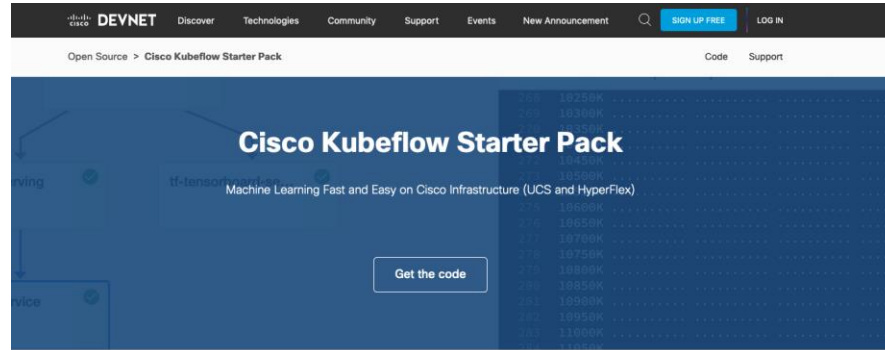
Username

Password

Connect

[Contacts](#) | [Feedback](#) | [Help](#) | [Site Map](#) | [Terms & Conditions](#) | [Privacy Statement](#) | [Cookie Policy](#) | [Trademarks](#)

# The Cisco kubeflow Starter Pack



Operationalize your machine learning workflows faster with  
Cisco Kubeflow starter pack



#### Manage Machine Learning at Scale

- Scalable deployment of ML pipelines
- Multi-user isolated environments
- Ability to control and roll back unwanted changes



#### Expedite ML from Idea to Production

- Consistent and reproducible experiments
- Better coordination between multiple team
- Efficient and easy way to deploy your models in production

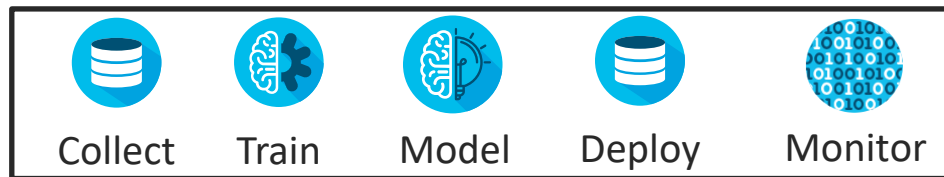
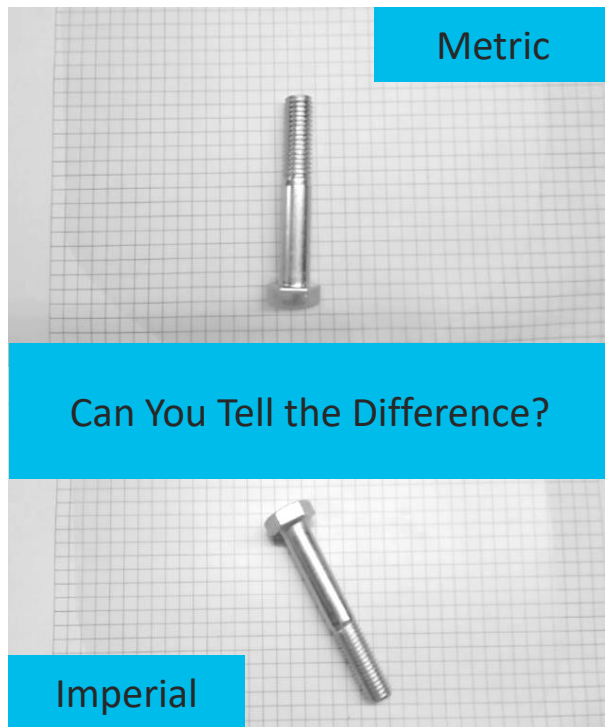


#### Data Pipeline All the Way to the Edge

- Training and inferencing closer to the data source
- Lower latency and reduced bandwidth consumption
- Server policies from the cloud with Cisco infrastructure

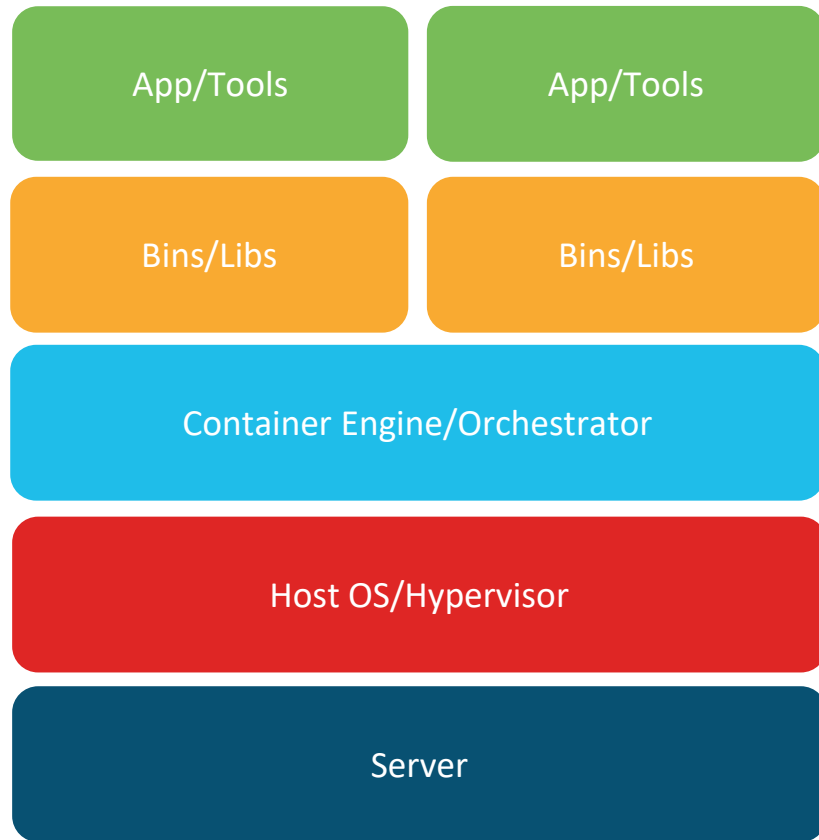
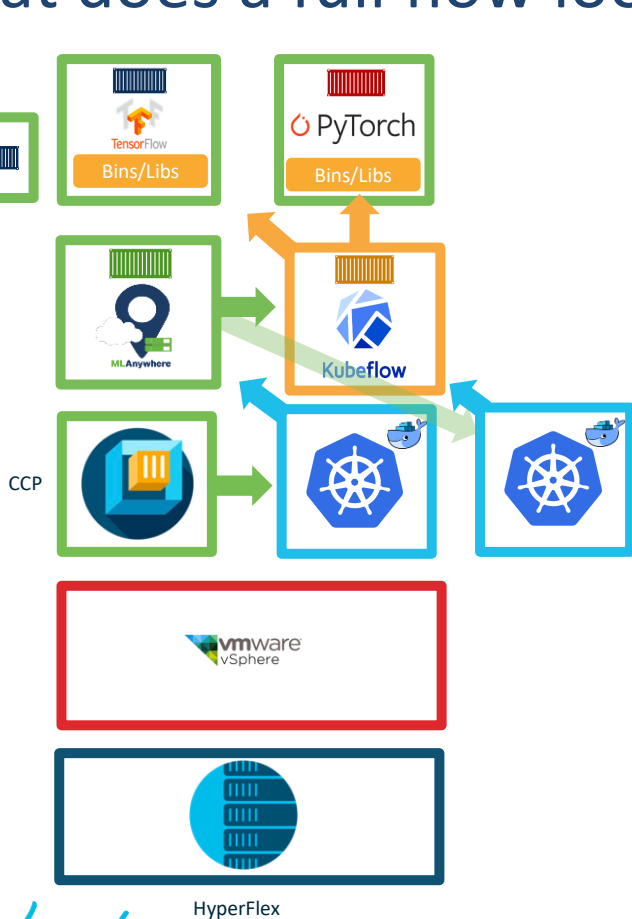
Preview the Cisco Kubeflow Starter Pack

# Recognizing Bolts Based on Imperial vs. Metric



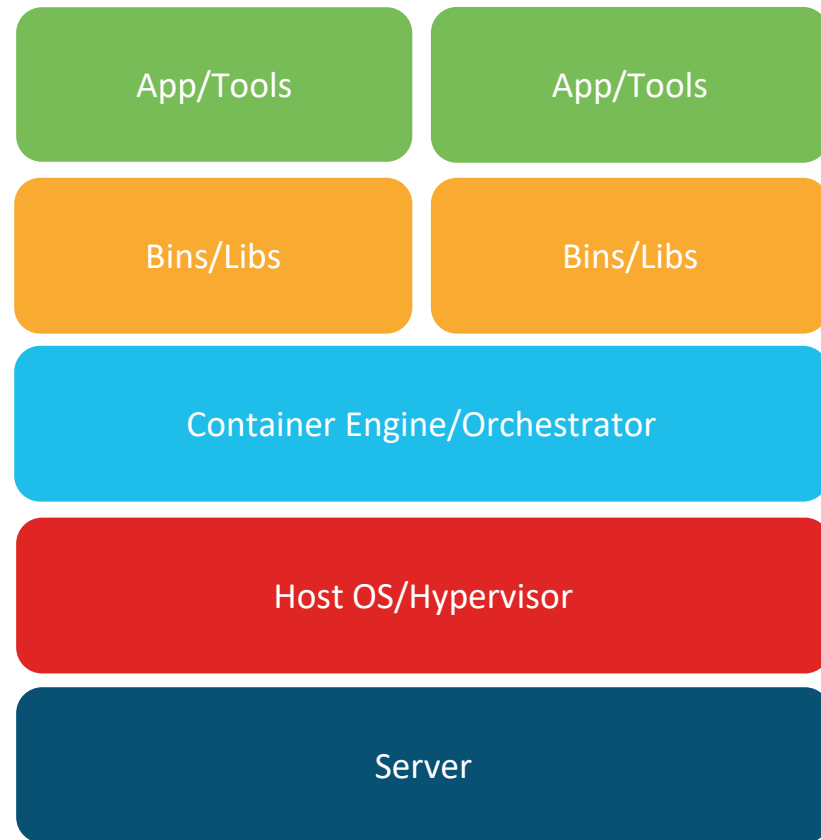
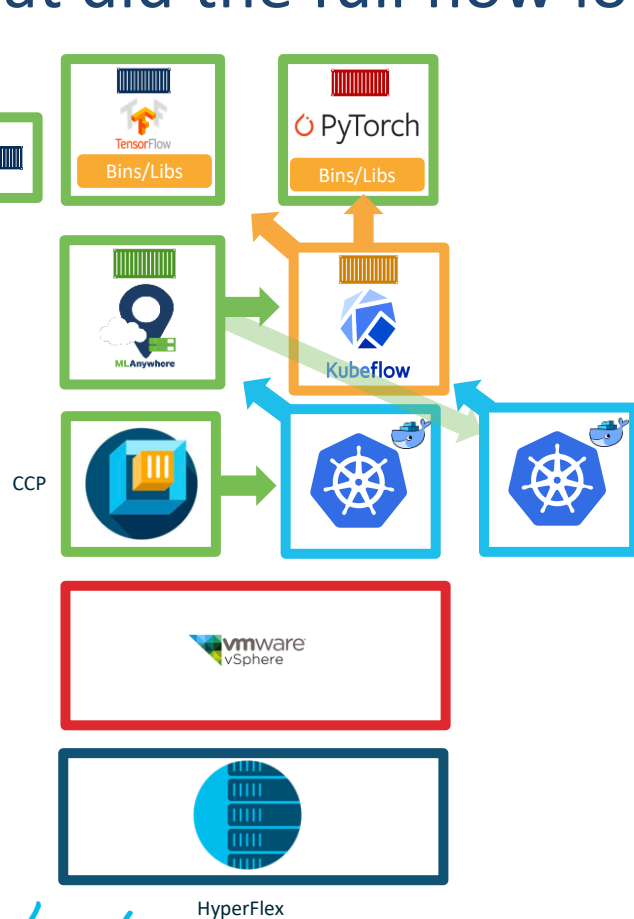
- Bolts based on (Imperial) inches vs. (Metric) centimeters are hard to distinguish: Wrong bolt can ruin equipment
- Use machine learning image classification to identify different types of bolts
- Kubeflow workflow for training, model evaluation, and inferencing
- Run on Cisco HX & CCP

# So what does a full flow look like?



Demo Time

# So what did the full flow look like?





# Cisco ML Solutions

# Cisco Container Platform



Turnkey Solution  
for Production-Grade Container  
Environments

## Native Kubernetes (100% Upstream)

Direct updates and best practices from open source community

## Multicloud Optimized

Deploy on-premise and native Amazon EKS, Azure AKS, Google GKE clusters

## Integrated

Networking | Storage | Management | Registry | Security | AI/ML

## Flexible Deployment Model

VMware (air gapped) | OpenStack (air gapped) | Public cloud

Easy to acquire, deploy and manage | Open and consistent | Extensible platform | World-class advisory and support

# Cisco Container Platform Feature Set

## Kubernetes-as-a-Service



### Setup

- Deploy Kubernetes clusters on HyperFlex, vSphere, OpenStack (CVIM), EKS, AKS
- CNI and Istio service mesh
- Persistent storage
- L4 / L7 Load Balancing
- Container Registry



### Consume

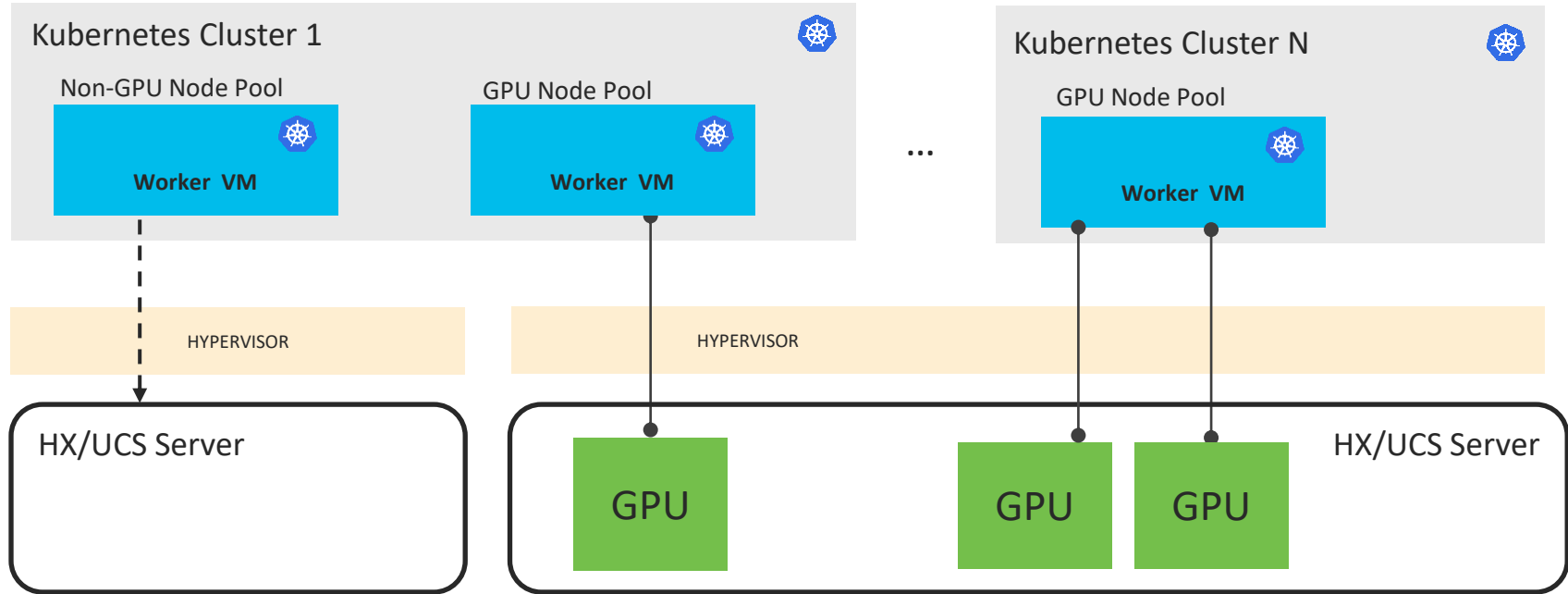
- AD Authentication / RBAC
- Resource based node pools
- Multi-GPU –as-a-Service
- Kubeflow (tech preview)
- UI – Kubernetes, API
- Security (policies, encryption)



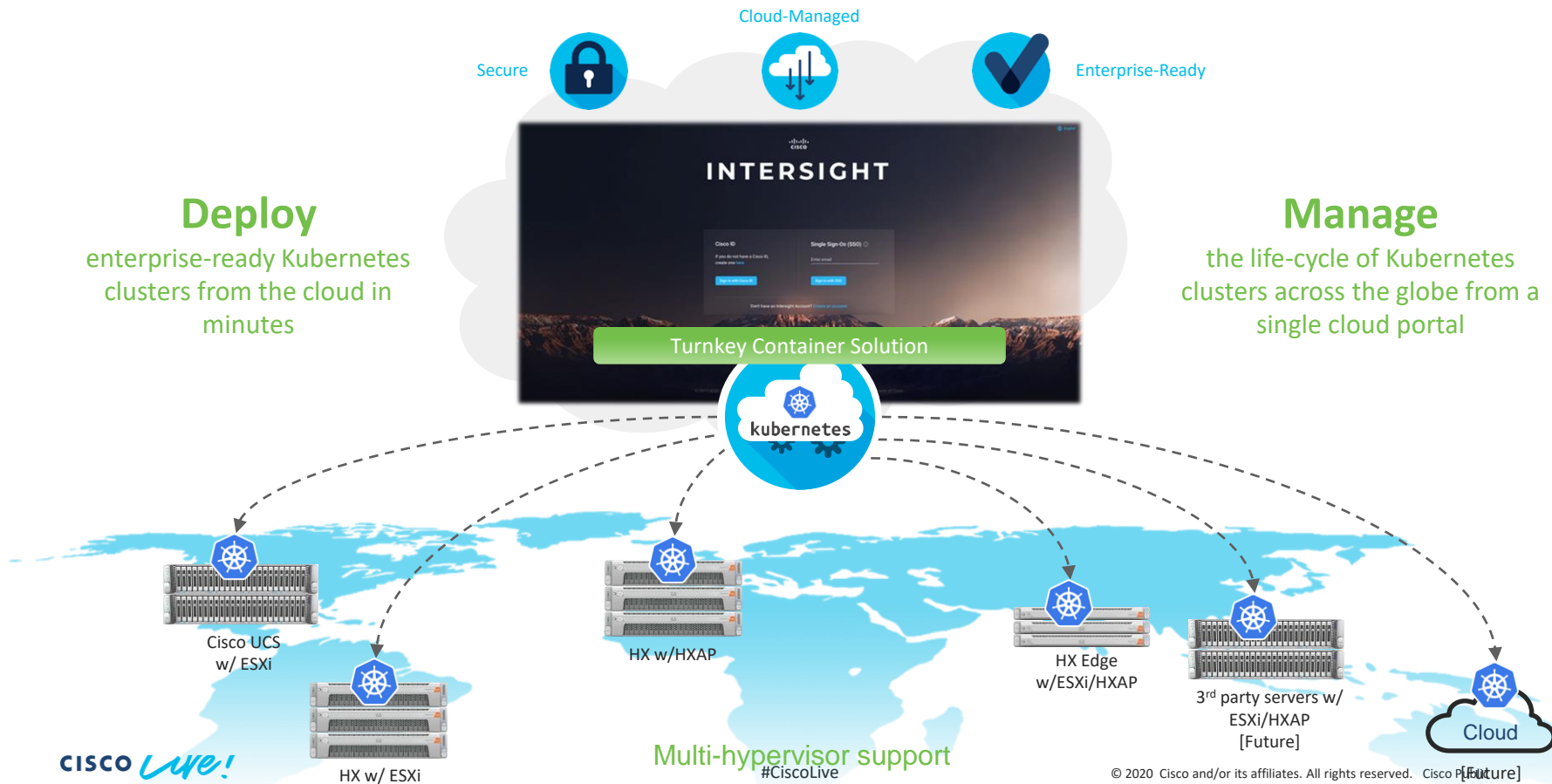
### Manage

- Add / remove Kubernetes nodes
- Lifecycle management (OS updates, Kubernetes upgrades)
- Prometheus/Grafana Monitoring
- EFK Logging
- Self-healing Kubernetes clusters
- Multi-master nodes

# Multi-GPU as a Service



# Cisco Intersight Kubernetes Service (formerly CCP)



# UCS AI/ML Portfolio Compute:



UCSM and Intersight Managed

Cisco  
IMC

XML  
API



## Test & Dev and Model Training

C240



HyperFlex 240

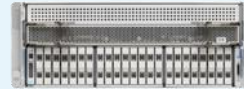


## Deep Learning/ Training

C480



C480 ML



## Inferencing

C/HX 220  
C/HX 240



Validated AI/ML SW For Turnkey (Working with Partners)



Better Together, Customer Choice, Cisco Validated Design with Eco-system





Opportunities to look out  
for



# Consistent Cloud and On-Premise Machine Learning Experience



Cloud

- Fast deployment
- Test-dev
- Simplicity
- APIs
- One time

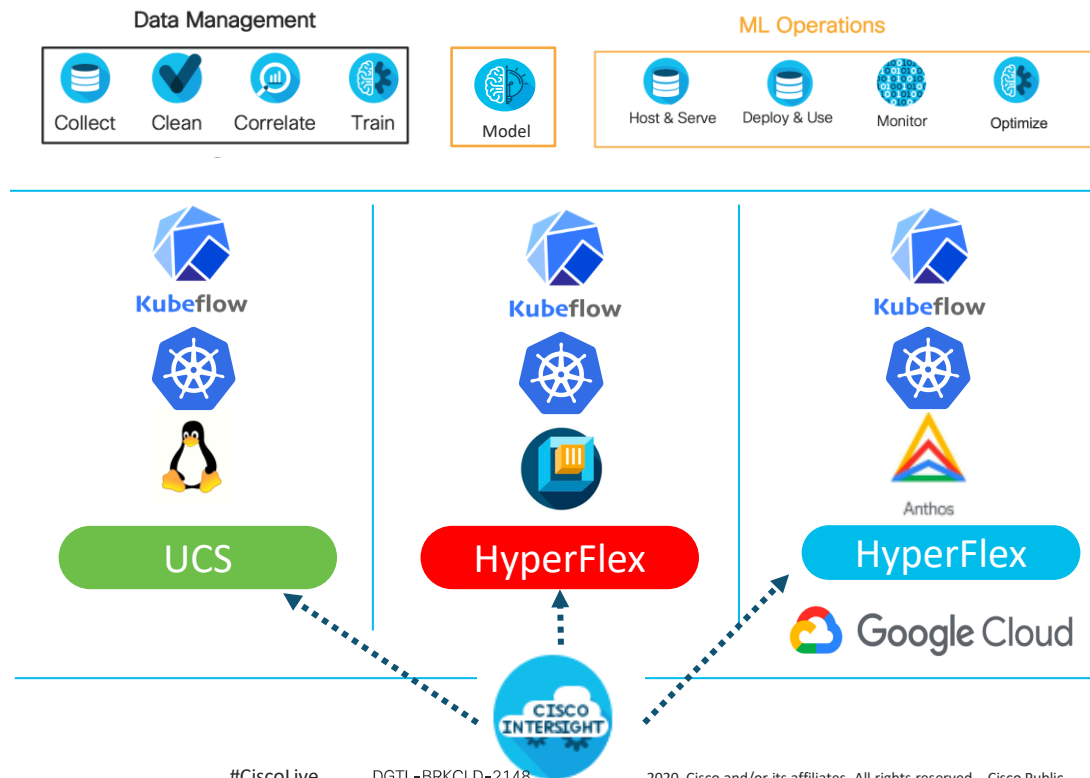


On Prem

- Data gravity and integration
- TCO and Performance
- Governance
- Remodel, retraining
- At scale production

CISCO *Live!*

## Cisco Hybrid Cloud Options



# Dedicated vs. Shared Infrastructure



Dedicated

Shared

## Data Scientists Experience

- Dedicated workstation vs. Share ML infrastructure

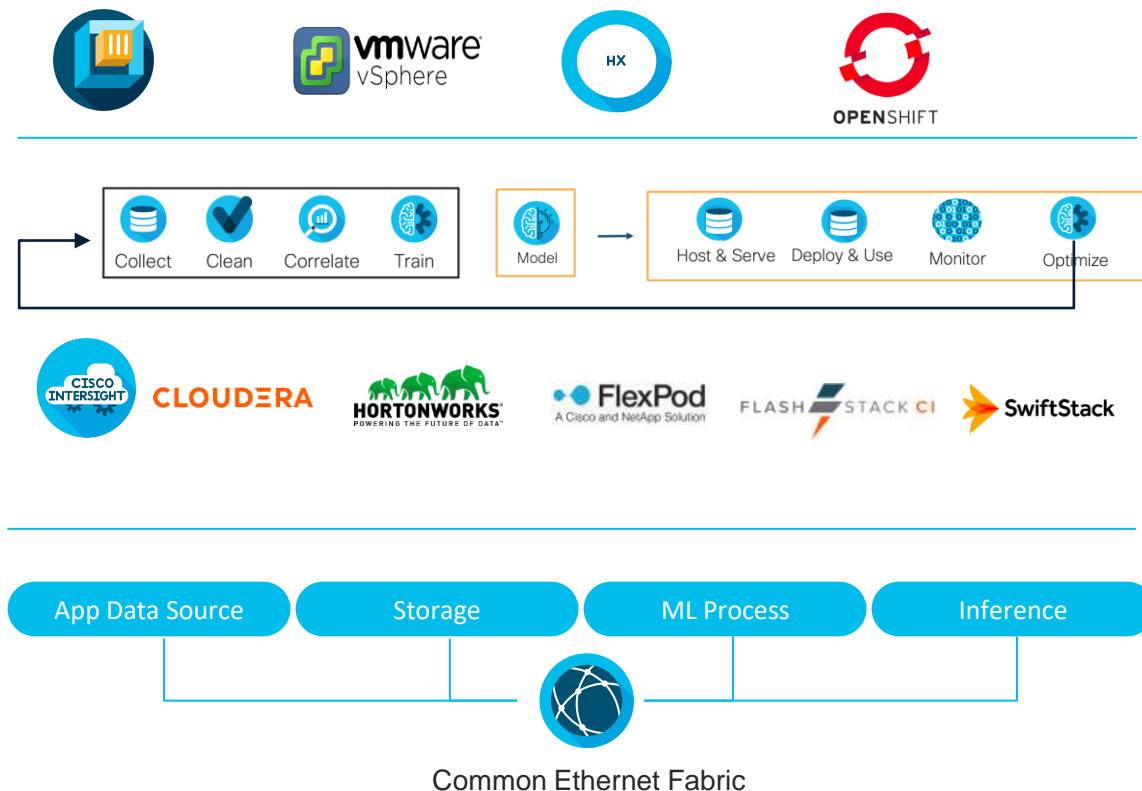
## Data Center Architecture

- Silo ML infra vs. Integrated Data Center Architecture

## Networking Architecture

- InfiniBand vs. Scalable Ethernet

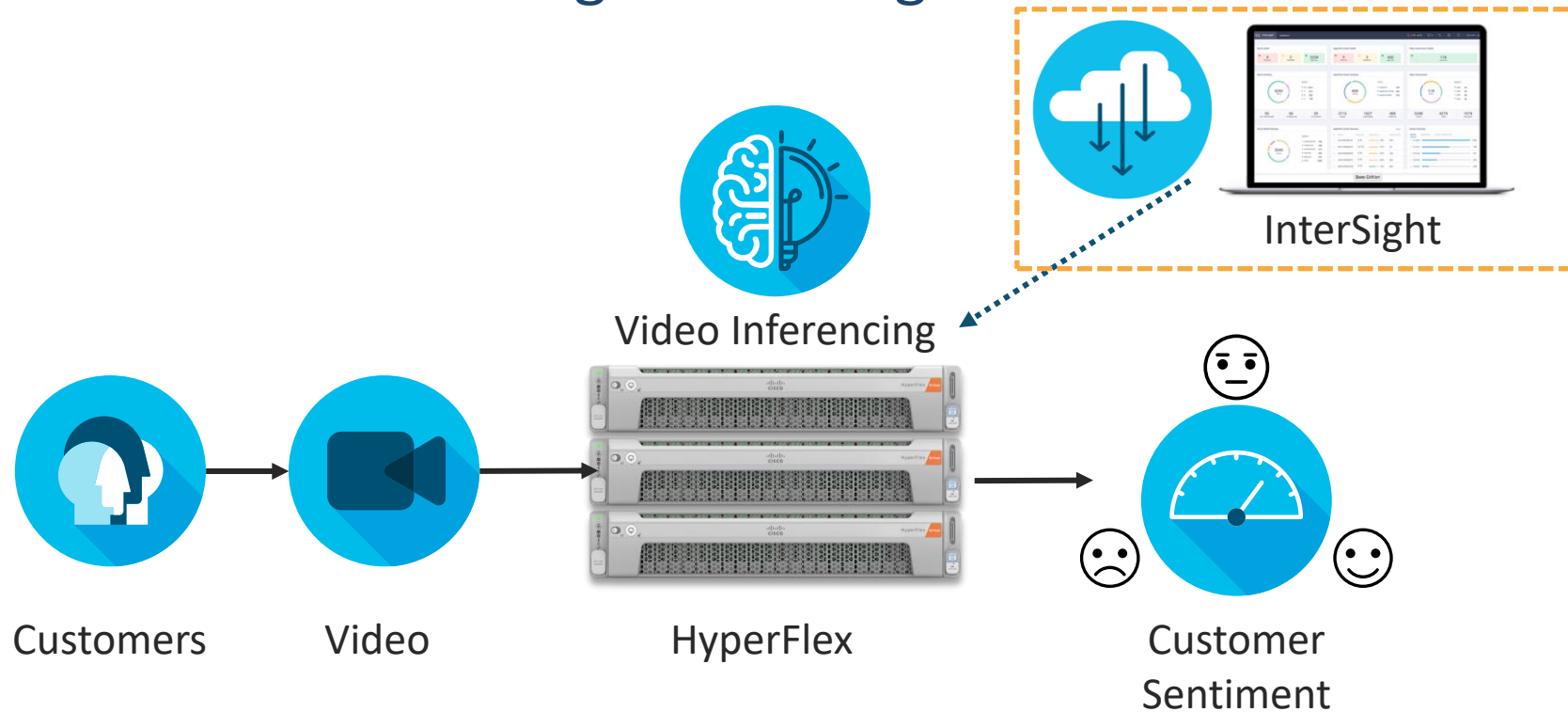
## Cisco: Unified Architecture



The background is a dark blue field filled with numerous small, semi-transparent squares and dots in various colors including light blue, green, yellow, orange, and red. These elements are scattered across the frame, with a higher concentration of yellow and orange squares forming a diagonal streak from the top right towards the bottom right.

The new frontier...The  
Edge!

# HyperFlex for Inferencing on the Edge



# Resources

- Cisco AI/ML Page  
<https://www.cisco.com/c/en/us/solutions/data-center/artificial-intelligence-machine-learning/index.html#~resources>
- Recent NVIDIA blog re: a100  
<https://blogs.cisco.com/datacenter/nvidia-a100-gpu-maximizes-ai-infrastructure-for-cisco-customers>
- Enterprise Study: Accelerating the AI Journey with Cisco  
<https://www.cisco.com/c/dam/en/us/solutions/collateral/data-center-virtualization/unified-computing/esg-wp-cisco-ai-ml-ucs.pdf>

Thank you



# Possibilities

#CiscoLive