# Bring Your Own Models (BYOM)- Machine Learning as a Service

**Malini Bhandaru, Soila Kavulya & Luis Daniel Castellanos**
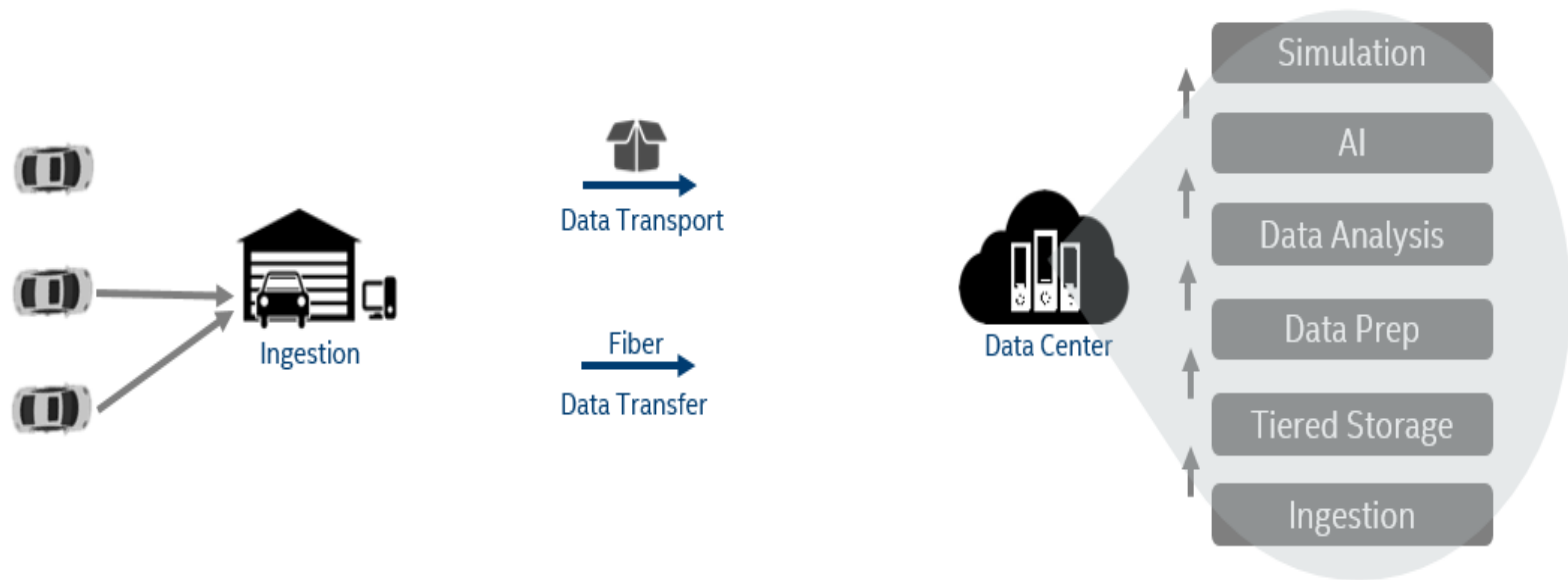**Contributors: Konrad Kurdej, Weiting Chen**
**INTEL**

**#ML9SAIS**

# Machine Learning Everywhere!



- Autonomous Vehicles

- Genomics

- Finance

- Supply Chain

# Autonomous Vehicles R&D Data Center

# Big Data

1 – 20 TB/car/hour
- Weather Conditions
- Time of Day
- Road Conditions
- Location
- Edge Cases

Object Detection Models
Environment Models
Driver Models
Privacy Preservation
Models

Image credit: https://www.wowwoodys.com/our-future-autonomous-cars/



**Under the bonnet**
How a self-driving car works

Signals from **GPS (global positioning system)** satellites are combined with readings from tachometers, altimeters and gyroscopes to provide more accurate positioning than is possible with GPS alone

**Lidar (light detection and ranging)** sensors bounce pulses of light off the surroundings. These are analysed to identify lane markings and the edges of roads

**Radar sensor**

**Video cameras** detect traffic lights, read road signs, keep track of the position of other vehicles and look out for pedestrians and obstacles on the road

**Ultrasonic sensors** may be used to measure the position of objects very close to the vehicle, such as curbs and other vehicles when parking

The information from all of the sensors is analysed by a **central computer** that manipulates the steering, accelerator and brakes. Its software must understand the rules of the road, both formal and informal
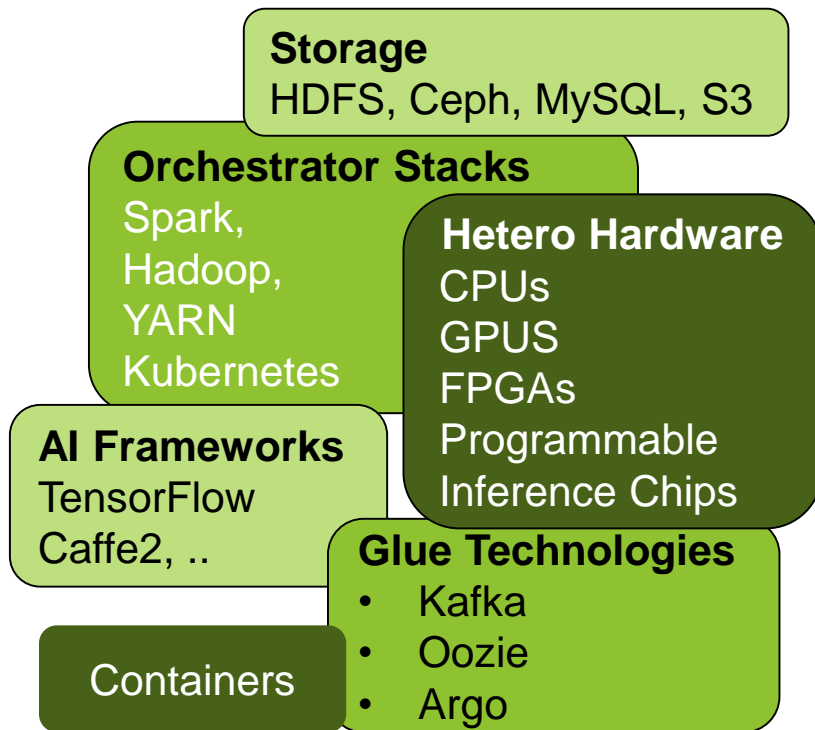
**Radar sensors** monitor the position of other vehicles nearby. Such sensors are already used in adaptive cruise-control systems
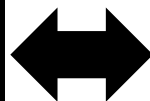
Source: *The Economist*

SPARK+AI SUMMIT 2018

#ML9SAIS

# Data Center Platform

- Fungible, Dynamic, Fast,
- Resilient
- Easy to Use

**Storage**
HDFS, Ceph, MySQL, S3

**Orchestrator Stacks**
Spark,
Hadoop,
YARN
Kubernetes

**Hetero Hardware**
CPUs
GPUS
FPGAs
Programmable
Inference Chips

**AI Frameworks**
TensorFlow
Caffe2, ..

**Glue Technologies**
- Kafka
- Oozie
- Argo

Containers

# Models Galore, Usages Rich

Input data → Model-1 → { Model-2 / Model-3 } → Simulation → Storage
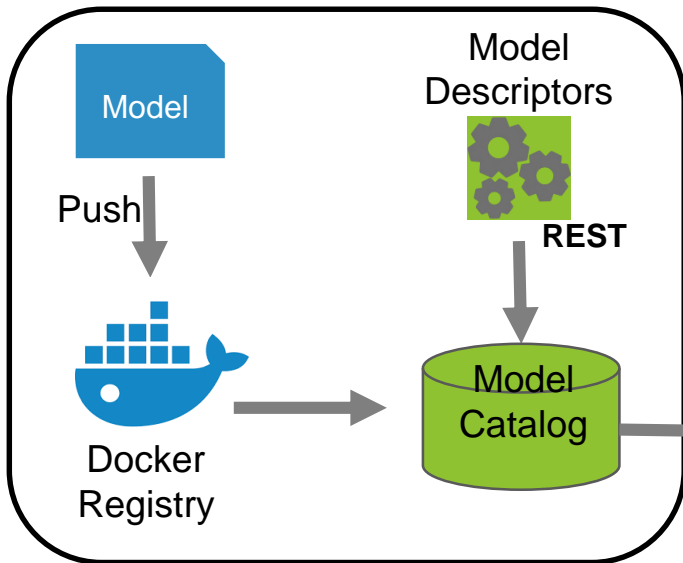


**MODEL**
**Name, Description,**
**type (mega | lean)**
**Framework & Version**
**Input, Output**
**ImageID: Container Registry ID**
**Training_Sets: { S1, S2 , S3}**
**Training_Label_Freq {L1:f1, L2:f2 ..}**
**Validation_Sets: {V1, V2}**
**Accuracy, Recall, Precision,**
**Speed, Size**
**Infrastructure:CPU/GPU/FPGA ..**

SPARK+AI
SUMMIT 2018

# Resources & API

**Model**
- **CRUD, Validate dependencies,**
- name, description, framework, version, hardware preference
- Tags (sharable, input-sensor ..)

**Data Transformer**
- **CRUD**
- Image resizer, compression, crypto, ..

**Pipeline**
- **CRUD, Start/Stop/Pause**
- workflow specifications, language

**Dataset**
- **CRUD**
- Name, description, data location
- (s3,hdfs, local file system)

Apache Beam, Spark, Argo

# Model Deployment Pipeline

**Model Creation and Registration**



Model

Push

Docker Registry

Model Descriptors

**REST**

Model Catalog

**Batch/Real-time Inference**

Sensor Data

Pre-processing
(e.g., RosBag parsing,
Video decode)

**Searchable object tags**

Inference

Post-processing

Storage

Web UI

# BYOM Options: Monolithic

**Spark Service**

Pod
Container
- Spark Driver
- ML Framework
- Model

Pod
Container
- Spark Executor
- ML Framework
- Model

Pod
Container
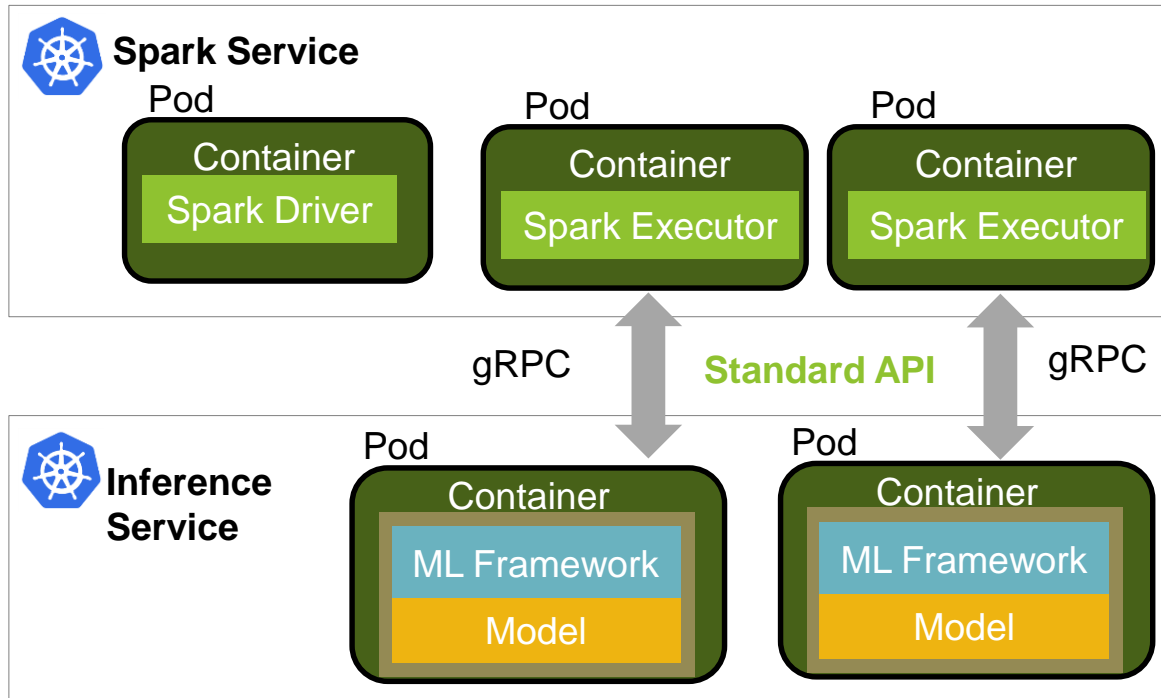- Spark Executor
- ML Framework
- Model

**Pros**
- Simple deployment
- Container life-cycle in full sync with workload
- No version tracking or mismatch concerns
- Data locality

**Cons**
- Larger container footprint
- Tight coupling between model and Spark engine

# BYOM Options: Just-In-Time Compose

**Spark Service**

Pod
Container
Spark Driver

Pod
Container
Spark Executor

Pod
Container
Spark Executor

gRPC  **Standard API**  gRPC

**Inference Service**

Pod
Container
ML Framework
Model

Pod
Container
ML Framework
Model

**Pros**
- Small container footprint
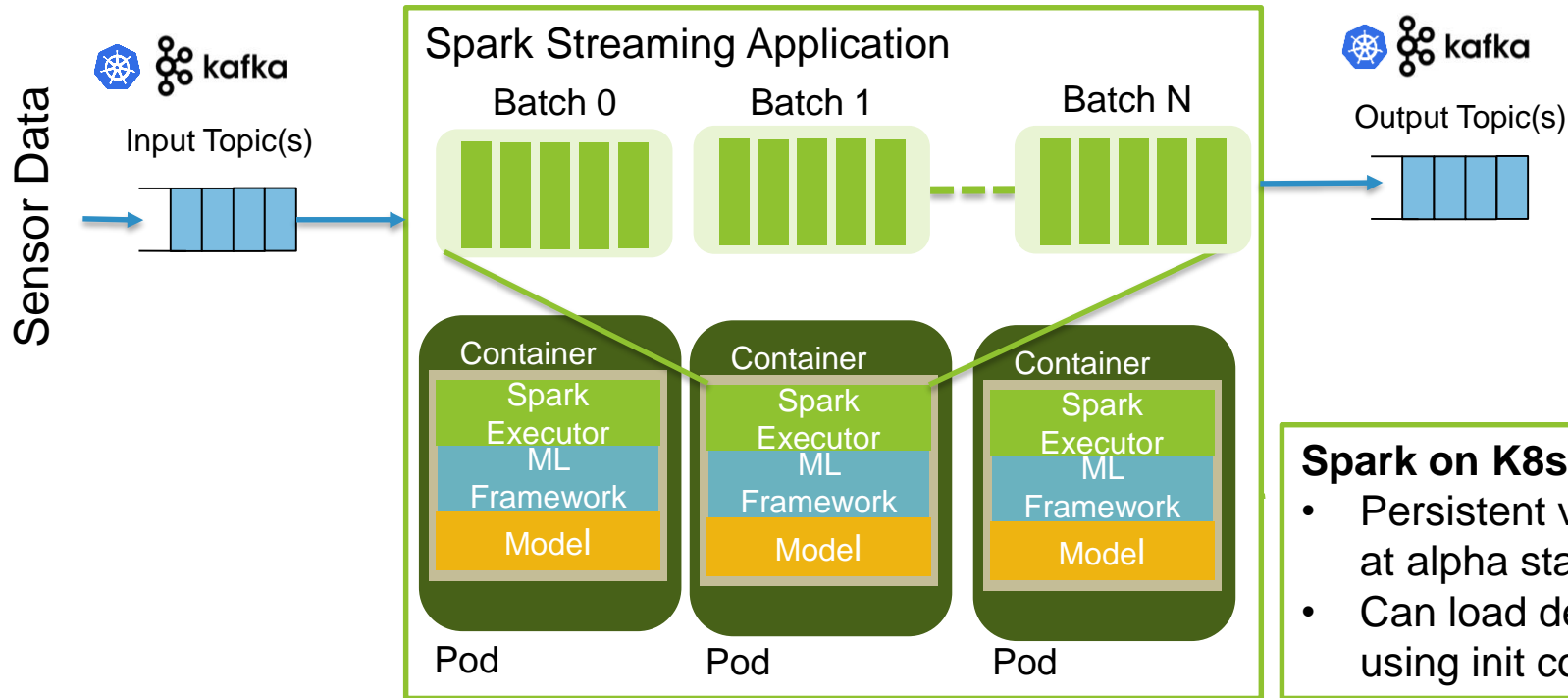- Multi-framework friendly
- Auto scales
- Standard API ⭐

**Cons**
- More complex orchestration workflow
- Additional mechanisms needed for data locality, e.g., pod affinity

# Deployment on K8s

# Demo

https://videoportal.intel.com/media/0_70vbt74e

# Future

- Support for Just-in-time-Composition
  - Tackling dependencies
- Resource scheduling, HW accelerator aware
- Hardware specific models (CPU/GPU/FPGA ..)
- Pipeline options with speed, accuracy, and resource availability projections

# Conclusion

- Across domains Bring-your-own-model genuine need for both R&D and Production Systems

- System and Orchestration developers typically not Machine Learning specialists –

  - reduce the barrier to adoption

# Thank You!

Kavulya, Soila P soila.p.kavulya@intel.com

Luis Daniel Castellanos (luis.daniel.castellanos@intel.com)

Kurdej, Konrad konrad.kurdej@intel.com

Bhandaru, Malini K malini.k.bhandaru@intel.com

Chen, Weiting <weiting.chen@intel.com>

# Please join us in the BYOM effort!

# References

- Autonomous Driving : https://medium.com/@andrewng/self-driving-cars-are-here-aea1752b1ad0, https://www.wired.com/story/embark-self-driving-truck-deliveries

- Docker Registry: https://docs.docker.com/registry

- Spark on Kubernetes: https://github.com/apache-spark-on-k8s/spark

- HDFS on K8: https://spark-summit.org/2017/events/hdfs-on-kubernetes-lessons-learned/

- RDD: https://www.slideshare.net/datamantra/anatomy-of-rdd

- TensorFlow Serving: https://www.tensorflow.org/serving/

- https://jaceklaskowski.gitbooks.io/mastering-apache-spark/content/spark-dynamic-allocation.html

- Kubeflow: ML toolkit for Kubernetes: https://github.com/kubeflow/kubeflow, https://cloud.google.com/blog/big-data/2018/03/simplifying-machine-learning-on-open-hybrid-clouds-with-kubeflow

- https://beam.apache.org/documentation/pipelines/design-your-pipeline/

# Upstreamed: Dynamic Resource Allocation

## Weiting Chen

*Resources are allocated at start but applications can request change at runtime.*
*Dynamic resource allocation uses shuffle service container for data shuffle (instead of Docker storage)*



**Kubelet**

App1 Driver Pod

*Step2*

*Step3*

Spark-Submit App1

*Step1*

**Kubernetes Master**

*Step4*

**Kubelet**

App1 Executor Pod

Shuffle Service Pod

*Step4*

**Kubelet**

App1 Executor Pod

Shuffle Service Pod

*Step4*

**Kubelet**

App1 Executor Pod

Shuffle Service Pod

*Design: Shuffle Service long running daemon pod.*
*Enable dynamicAllocation, specify min, max, and*
*initial number of executors.*