

**RESILIENCE  
REALIZED**



KubeCon



CloudNativeCon

North America 2021



KubeCon



CloudNativeCon

North America 2021

RESILIENCE  
REALIZED

# AIoT Ops - Using Kubernetes and ML Ops to build Edge ML applications

*Asheesh Goja, Cisco*

# About Me

Senior Solutions Architect @ Cisco

Team : Emerging Technologies and Incubation



# Topics

## Talk – *Why, What and how of AIoT ML Ops*

- Why
  - AIoT
  - ML Ops
- What
  - Emergent behaviors
  - AIoT solution Patterns
  - Reference Architecture
- How
  - Reference implementation using K3S, TFLite and Argo
  - Deployment topology
  - Electrical circuit setup

## Demo - *Step by step guide to build and deploy an AIoT MLOps solution*

- Run the demo in a fully configured environment and show
  - ML workflow DAGs – Run Extract , Drift , Quantize and Train jobs
  - Perform Edge Inferences on Coral TPU cluster
  - Perform MCU inference on ESP32 RTOS
  - Introduce fault on the motor and detect it
- Tear it down the demo and show how to rebuild:
  - Infrastructure
    - Setup K3S and Strimzi
    - Setup Argo workflows
  - Platform
    - Setup model registry
    - Setup container registry
    - Build and deploy protocol bridge and FOTA
  - Application
    - Configure Argo DAGs to deploy Extract, drift detection, training and quantization modules
    - Build and deploy AIoT ML module to TPU and MCU devices
    - Show HA and Scaling capabilities

# AIoT – Artificial Intelligence of Things

## Problem – By 2025

- 55 billion IoT devices
- 73.1 Zettabytes of data
- Data deluge – more data less insight
- Privacy concerns
- Wide cyber attack vector with internet connected devices

## The solution - Embed ML in IoT

- Perform inferences close to sensors and actuators
- Minimize need for cloud connectivity
- Execute close loop anomaly detection and intervention
- Enable operations that require low latency communication
- Improve s/n ratio - Reduce inconsequential data uploads to the cloud
- Generate strong privacy protection – keep data on the device or edge tier



# Emergent behaviors and considerations

- Computational complexity impacts resource availability
  - Training vs. Inference complexity
  - Space, time and sample complexity
  - Computation, memory, bandwidth, energy
  - Accuracy, resource and complexity tradeoffs
- Limited memory, processing power and energy capacity impacts inference capability
- ML frameworks are too bulky for embedded devices
- Insufficient metrics to measure model performance
  - MACs and FLOPs inadequate
- Optimization strategies introduce training and model entropy
  - Quantization error
  - Errors introduced due to binarization of activations and weights
- Multiple incompatible computational hardware architectures
  - ARM64 vs x86 vs AMD64 vs ARM Cortex-M

| ML Model                   | Size            | Training                  | Inference              |
|----------------------------|-----------------|---------------------------|------------------------|
| Random Forest              | $O(N_{tree} m)$ | $O(N_{tree} mn (\log(m))$ | $O(N_{tree} (\log(m))$ |
| <b>Logistic Regression</b> | $O(n)$          | $O(mn^2 + n^3)$           | $O(n)$                 |
| Naive Bayes                | $O(nc)$         | $O(Mn + nc)$              | $O(nc)$                |
| Support-vector machine     | $O(n)$          | $O(m^2n)$                 | $O(m_{sv}n)$           |

*m samples, c classes and n dimensions*

## ML Ops

- Express ML Ops dependencies as DAGs
- Manage extract, drift detection, training, validation and quantization on the edge tier
- Use workflow pipelines for continuous evaluation and training
- Use hardware accelerator aware pod placement strategies

## Architecture

- Energy efficient hardware acceleration
- Containerization
- Streaming API sidecar
- Decouple embedded inference from communication
- Automated container orchestration on edge devices

## Embedded ML

- Decouple training from inference
- Introduce preprocessing Inferencing at analog tier
- Use early termination and cascaded models
- Binarize weights and activations at train time
- Use bit-shift operations over multiplications
- Explore model partitioning, edge caching and input filtering
- Perform data compression or sparsification
- Use Transfer Learning
- Use DSP to filter out noise

# AIoT ML Ops Reference Architecture

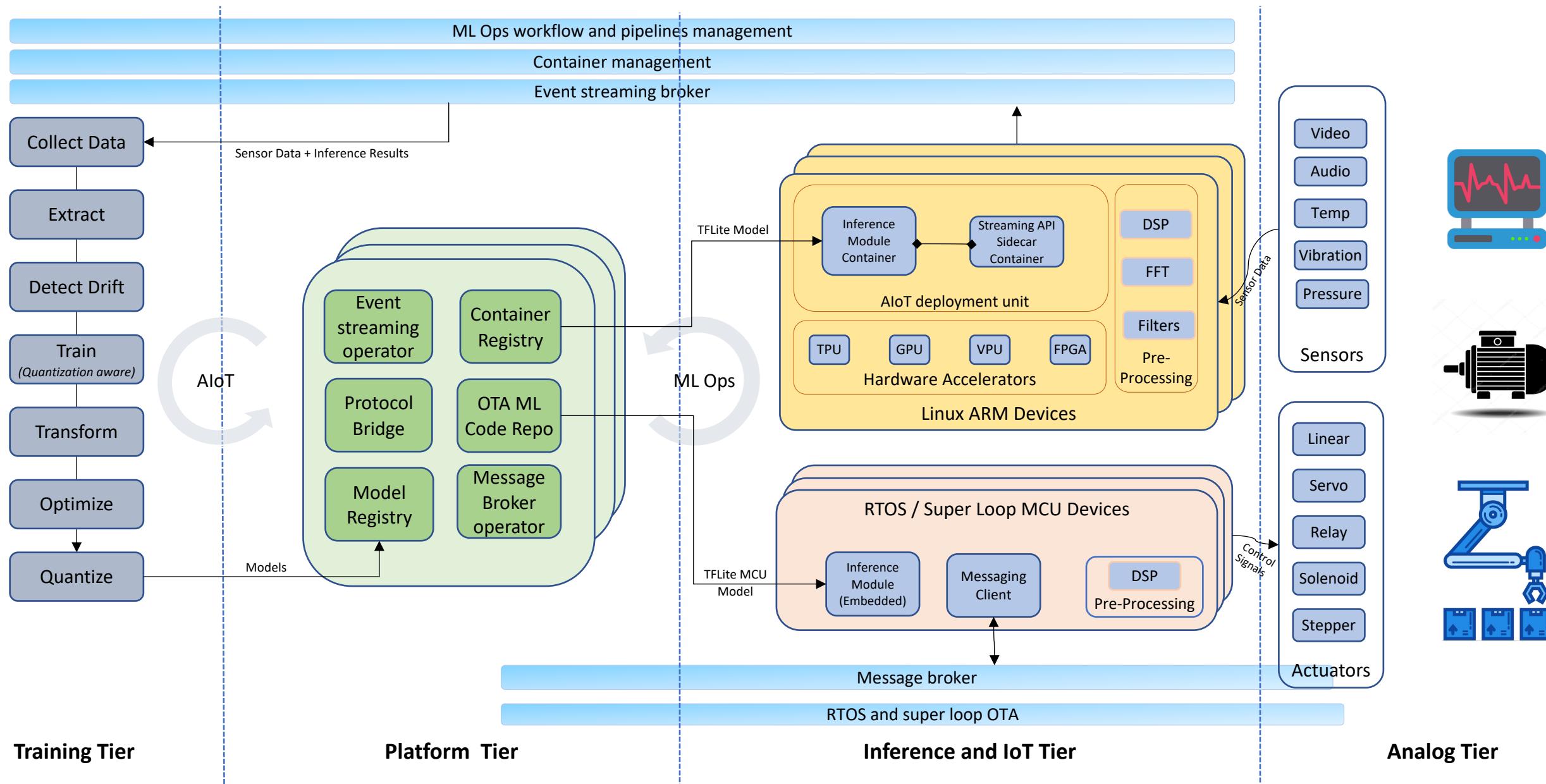


KubeCon



CloudNativeCon

North America 2021



# Reference Implementation- Strmizi, K3S, Argo and TFLite

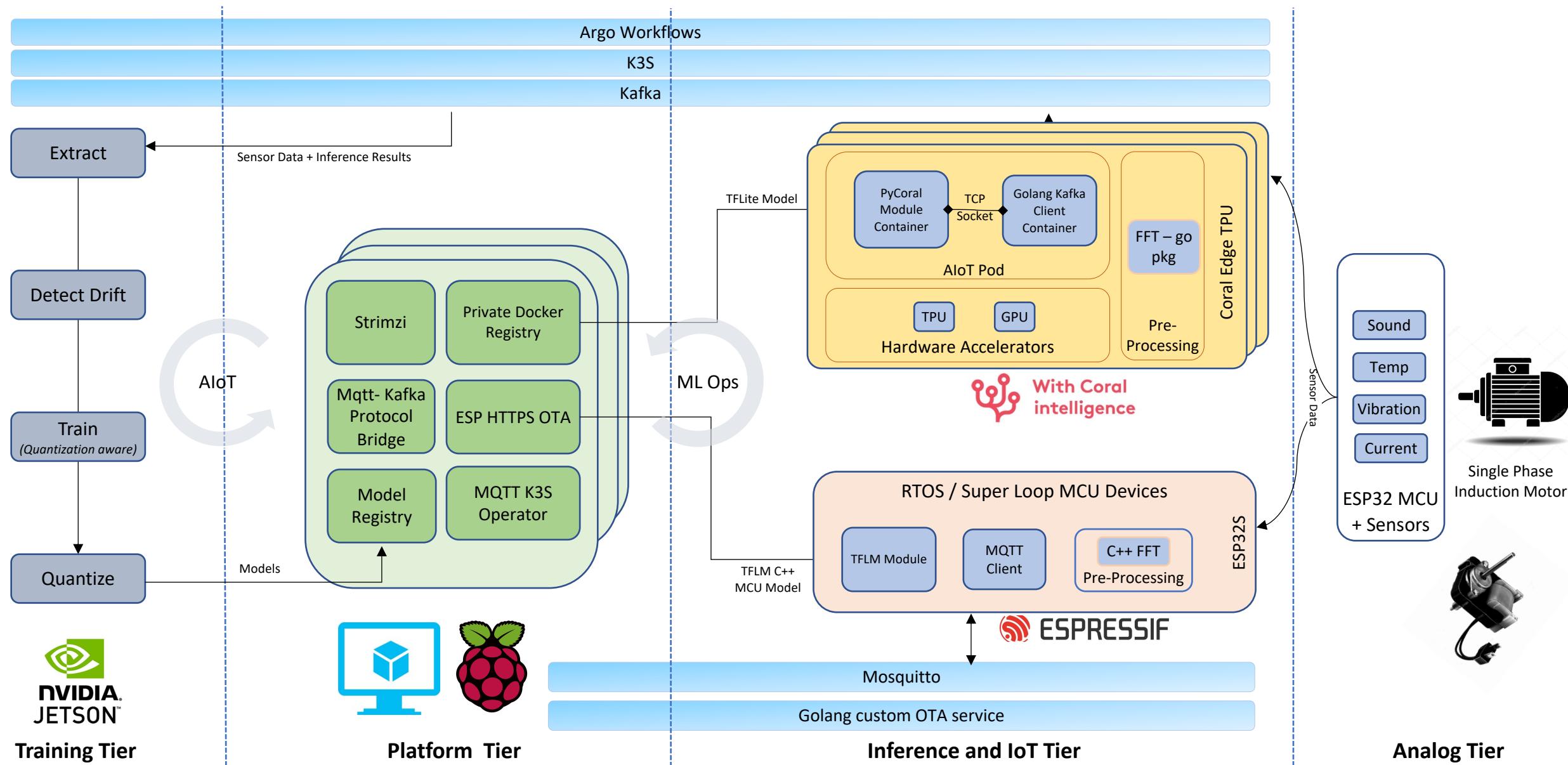


KubeCon



CloudNativeCon

North America 2021

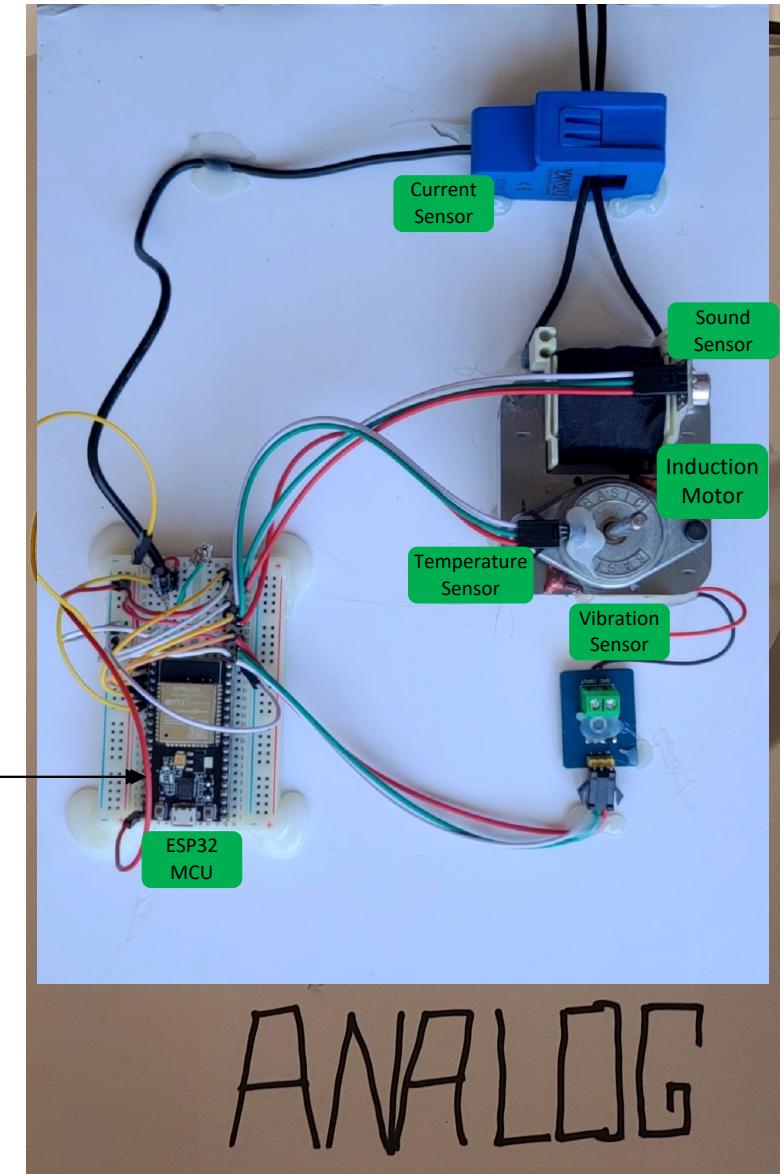
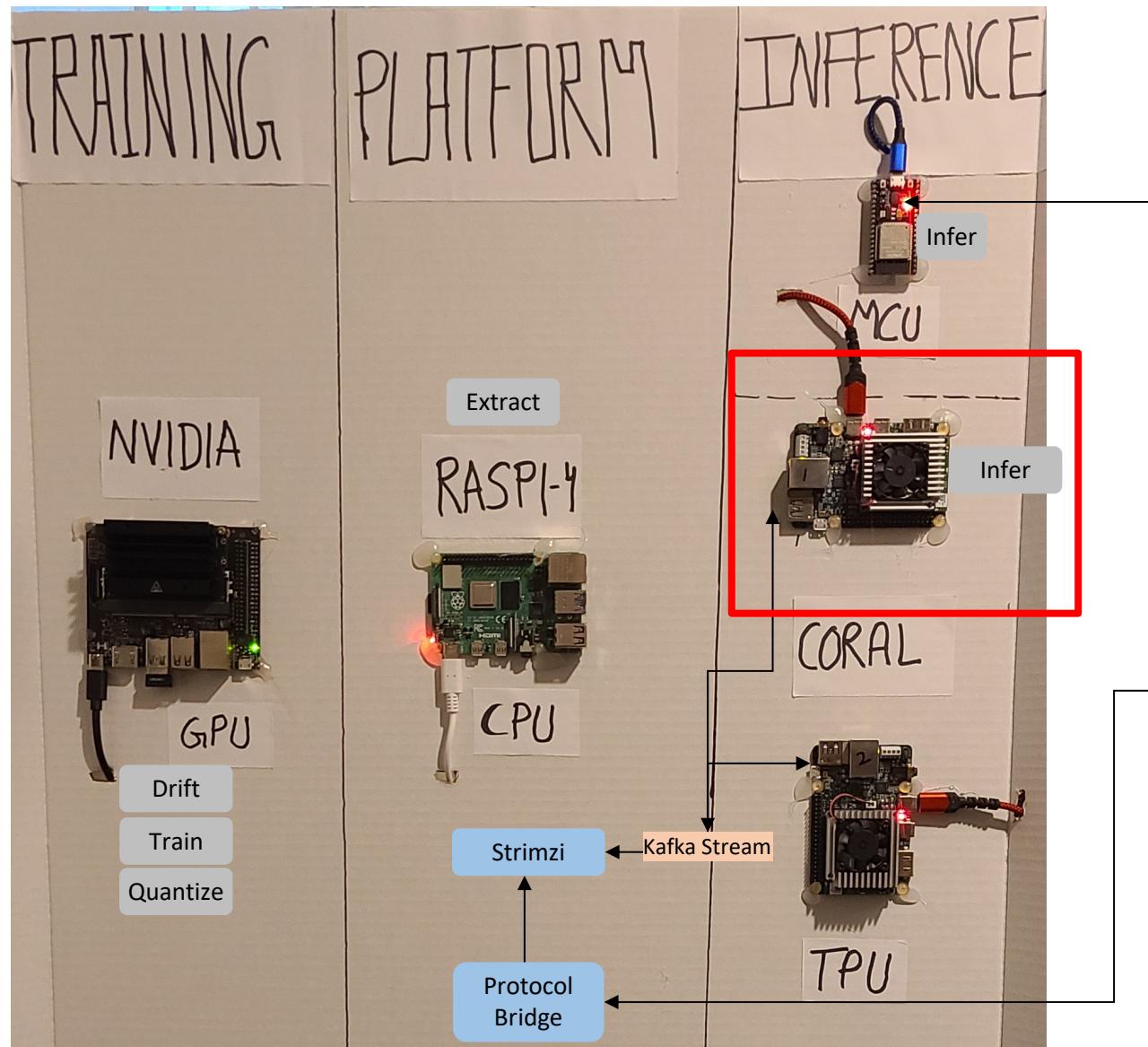


Let's start  
building...Demo Time

---



# Demo- On stage setup



# Under the hood – TPU and MCU Inference modules



North America 2021

```
#uncomment this section to te Quick Fix... (#)
from pycoral.adapters import classify
from pycoral.adapters import common
from pycoral.utils.dataset import read_label_file
from pycoral.utils.edgetpu import make_interpreter
from periphery import GPIO

ledB = GPIO("/dev/gpiochip2", 9, "out") # 16
ledG = GPIO("/dev/gpiochip4", 10, "out") # 18
ledR = GPIO("/dev/gpiochip4", 12, "out") # 22

try:
    current = np.float32(json.loads(msg)["current"])
    temperature = np.float32(json.loads(msg)["temperature"])
    vibration = np.float32(json.loads(msg)["vibration"])
    sound = np.float32(json.loads(msg)["sound"])

    telemetry = [current, temperature, vibration, sound]

    np_arr_64 = np.array(telemetry)
    np_arr_f32 = np_arr_64.astype(np.float32)

    inp_details = interpreter.get_input_details()
    out_details = interpreter.get_output_details()

    interpreter.set_tensor(inp_details[0]['index'], [telemetry])

    interpreter.invoke()

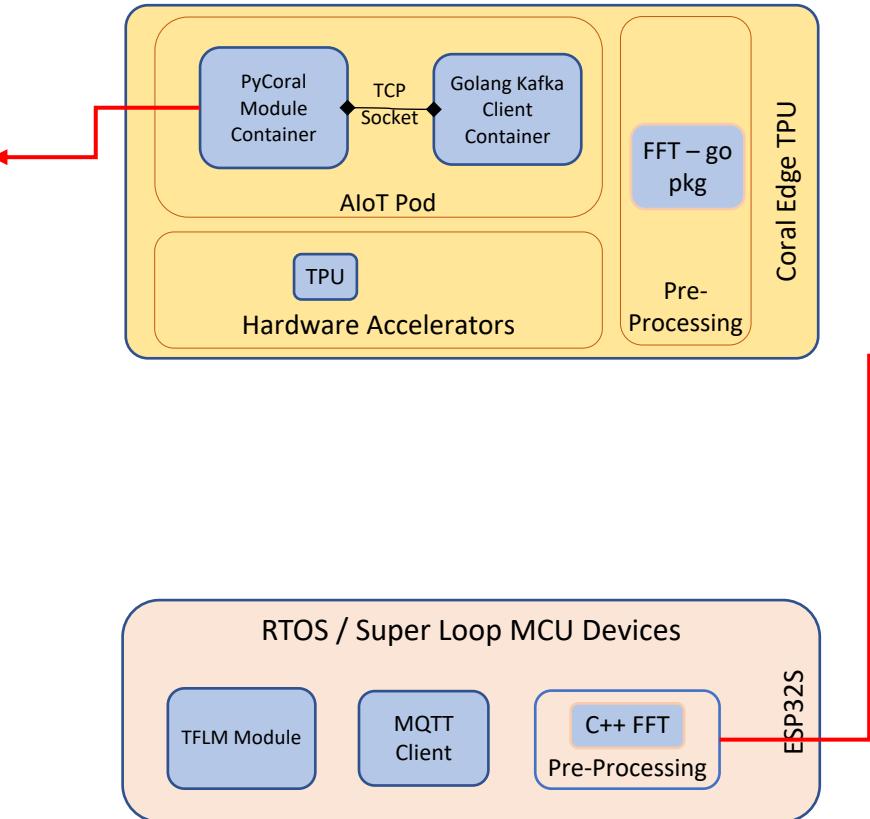
    output_details = interpreter.get_output_details()
    predictions = interpreter.get_tensor(output_details[0]['index'])

    output_index = interpreter.get_output_details()[0]["index"]
    ten = interpreter.get_tensor(output_index)

    fault = float('%.f' % ten)

    hn_ip_tuple = getHostnameAndIpAddress()

    logToFile(bcolors.BOLD +
              "Edge TPU logistic regression inference value : {}".format(fault))
```



PyCoral Logistic Regression Module

```
#include "tensorflow/lite/micro/all_ops_resolver.h"
#include "tensorflow/lite/micro/micro_error_reporter.h"
#include "tensorflow/lite/micro/micro_interpreter.h"
#include "tensorflow/lite/schema/schema_generated.h"
#include "tensorflow/lite/version.h"

interpreter = new tflite::MicroInterpreter( model, *resolver, tensor_a);

TfLiteStatus allocate_status = interpreter->AllocateTensors();
if (allocate_status != KTFLiteOk)
{
    TF_LITE_REPORT_ERROR(error_reporter, "AllocateTensors() failed");
    return;
}

size_t used_bytes = interpreter->arena_used_bytes();
TF_LITE_REPORT_ERROR(error_reporter, "Used bytes %d\n", used_bytes);

// Obtain pointers to the model's input and output tensors.
input = interpreter->input(0);
output = interpreter->output(0);

Serial.printf("mqtt message for TFLM inference - Current: %d, Temp: %d, Vibration: %d, S

input->data.f->getInputBuffer()[0] = current;
input->data.f->getInputBuffer()[1] = temperature;
input->data.f->getInputBuffer()[2] = vibration;
input->data.f->getInputBuffer()[3] = sound;

interpreter->Invoke();}
float result = output->data.f[0];

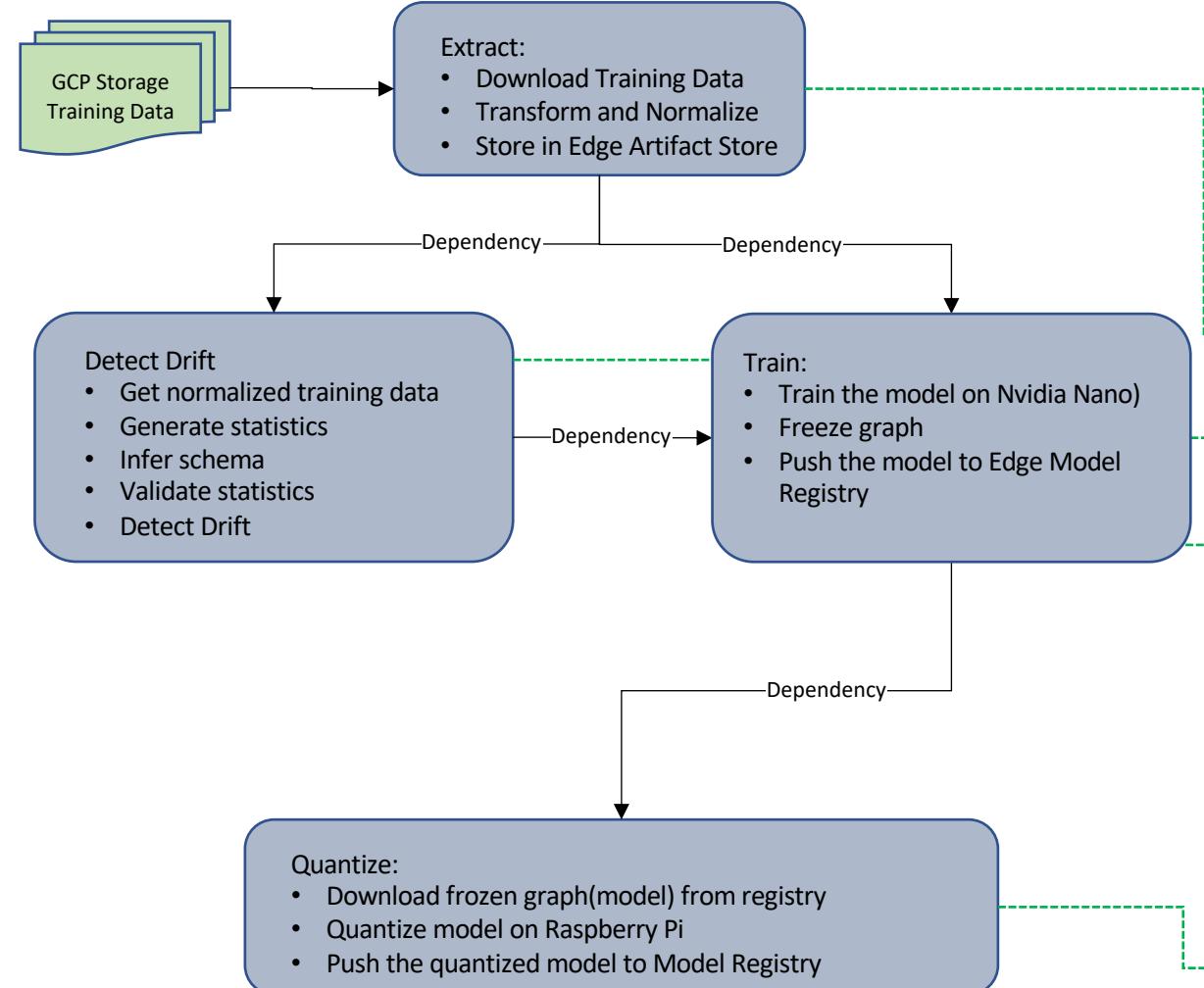
Serial.printf("performing tflm inference on the mqtt mesage result = %.2f \n", result);

if (result > 0.7)
{
    for (int i = 0; i < 30; i++)
    {
        digitalWrite(LED_BUILTIN, HIGH);
        delay(100);
        digitalWrite(LED_BUILTIN, LOW);
        delay(100);
    }

    char sensor_data[255];
    sprintf(sensor_data, 255, "{\"inference-level\": \"%d\", \"chipID\": \"%d\", \"current\": %d, \"chipId, current, temperature, vibration, sound");
    client.publish(outTopic, sensor_data);
```

C++ TFLM Logistic Regression

# Under the hood : Argo DAG - Train and Quantize Pipeline



```

1  apiVersion: argoproj.io/v1alpha1
2  kind: Workflow
3  metadata:
4    generateName: kubecon-aiotdemo-dag-
5  spec:
6    entrypoint: kubecon-aiotdemo-dag
7    templates:
8      - name: kubecon-aiotdemo-dag
9        dag:
10       tasks:
11         - name: extract
12           template: extract-template
13           arguments:
14             parameters:
15               - name: message
16               value: ""
17
18         - name: validate
19           dependencies: [extract]
20           template: validate-template
21           arguments:
22             parameters:
23               - name: message
24               value: ""
25
26         - name: train
27           dependencies: [validate,extract]
28           template: train-template
29           arguments:
30             parameters:
31               - name: message
32               value: ""
33
34         - name: quantize
35           dependencies: [train]
36           template: quantize-template
37           arguments:
38             parameters:
39               - name: message
40               value: ""
  
```

- name: extract-template  
 inputs:  
 parameters:  
 - name: message  
 container:  
 image: docker.35.202.84.160.nip.io:5000/extract\_module:latest  
 securityContext:  
 privileged: true  
 env:  
 - name: TRAINING\_DATA\_UPLOAD\_REGISTRY\_URL  
 value: "http://35.236.22.237:30007/uploadTrainingData"  
 - name: GCP\_BUCKET  
 value: "kubecon2021-aiot-mlops-demo"  
 nodeSelector:  
 kubernetes.io/hostname : "kubecon-aiot-control-node"

- name: validate-template  
 inputs:  
 parameters:  
 - name: message  
 container:  
 image: docker.35.202.84.160.nip.io:5000/validation\_module:latest  
 env:  
 - name: TRAINING\_DATA\_URL  
 value: "http://35.236.22.237:30007/training\_data"  
 nodeSelector:  
 kubernetes.io/hostname : "kubecon-aiot-control-node"

- name: train-template  
 inputs:  
 parameters:  
 - name: message  
 container:  
 image: docker.35.202.84.160.nip.io:5000/training-module:latest  
 env:  
 - name: MODEL\_REGISTRY\_URL  
 value: "http://35.236.22.237:30007/uploadModel"  
 nodeSelector:  
 kubernetes.io/hostname : "agentnode-nvidia-jetson"

- name: quantize-template  
 inputs:  
 parameters:  
 - name: message  
 container:  
 image: docker.35.202.84.160.nip.io:5000/quantize-module:latest  
 securityContext:  
 privileged: true  
 env:  
 - name: MODEL\_DOWNLOAD\_REGISTRY\_URL  
 value: "http://35.236.22.237:30007/full"

MLOps Train/Quantize Pipeline

Argo DAG YAML

# Edge ML Devices and Accelerators



KubeCon

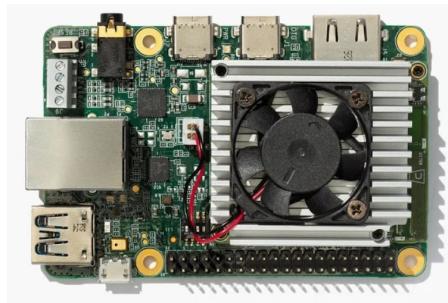


CloudNativeCon

North America 2021

## Google Coral Edge TPU

- CPU - Quad Cortex-A53 @ 1.5 GHz
- GPU - Vivante GC7000Lite
- TPU - Google Edge TPU coprocessor
- VPU - 4Kp60 HEVC/H.265
- Memory - 1 GB LPDDR4



## NVIDIA Jetson Nano

- CPU – Quad-core ARM® A57 @ 1.43 GHz
- GPU - 128-core NVIDIA Maxwell™
- Memory - 2 GB 64-bit LPDDR4



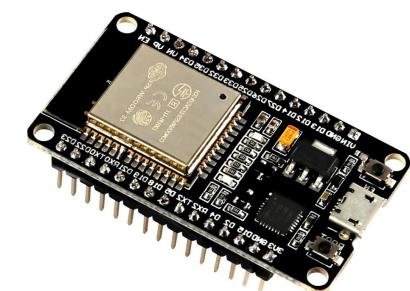
## Raspberry Pi 4

- CPU - Quad Cortex-A72 @ 1.5GHz
- Memory - 4GB LPDDR4



## ESP32S

- MCU – Dual Core Xtensa® 32-bit LX6 @ 40Mhz
- Memory - 448 KB ROM, 520 KB SRAM



# AIoT Demo – Deployment Topology

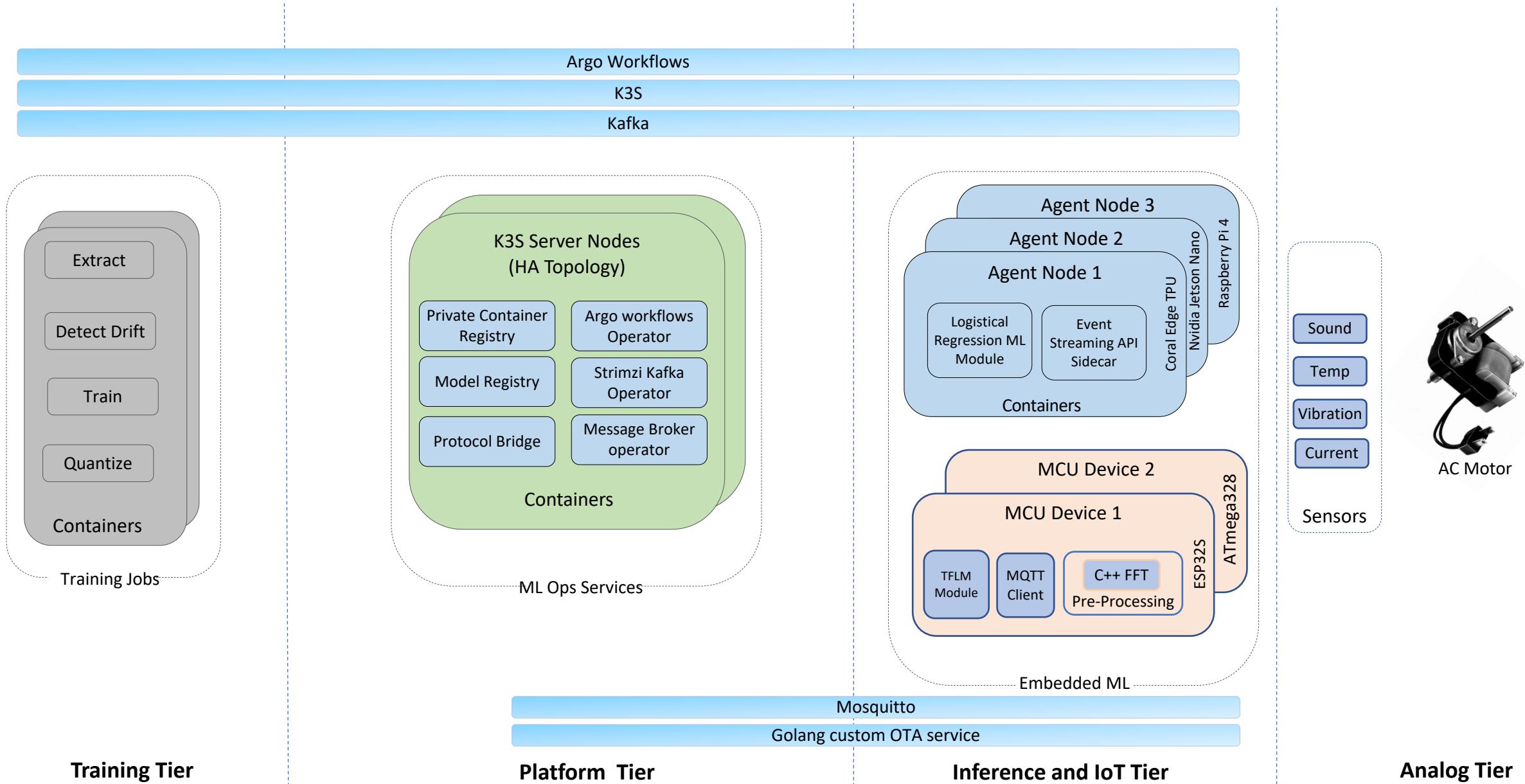


KubeCon

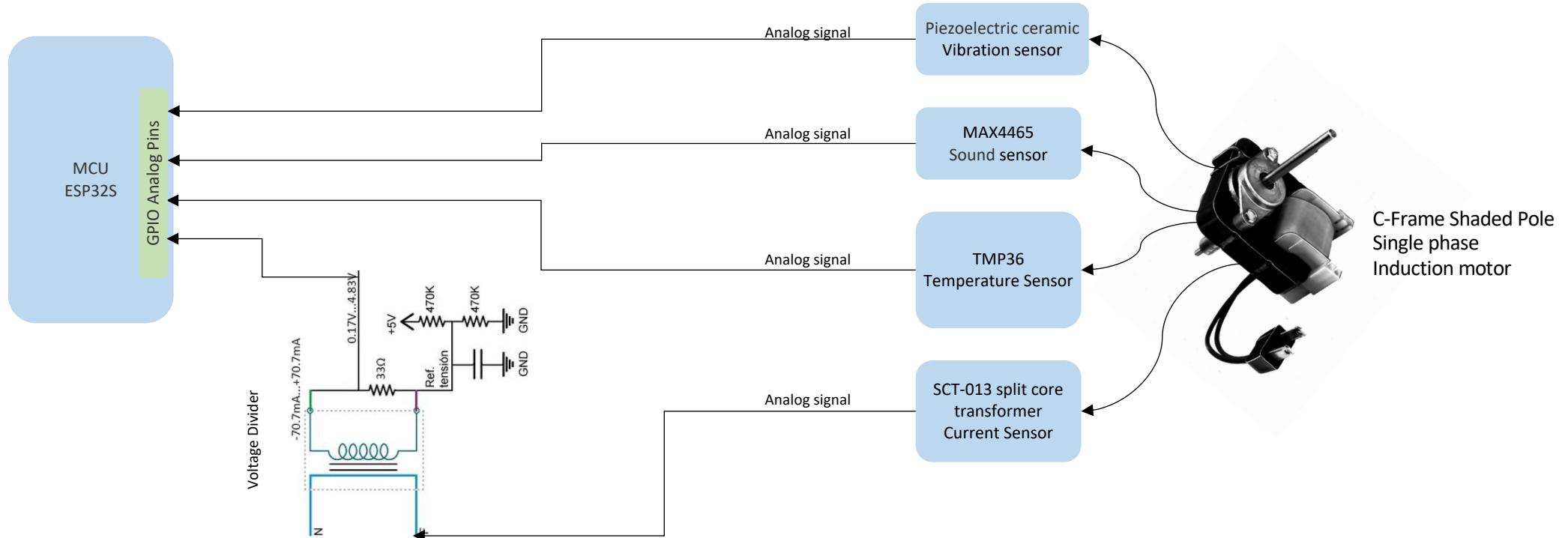


CloudNativeCon

North America 2021



# Electrical Circuit Setup



# Live Demo – K3S, Strimzi and Argo Setup



KubeCon



CloudNativeCon

North America 2021

## K3S Cluster

| NAME                      | STATUS | ROLES         | AGE   | VERSION      | INTERNAL-IP   | EXTERNAL-IP   | OS-IMAGE                       | KERNEL-VERSION  | CONTAINER-RUNTIME       |
|---------------------------|--------|---------------|-------|--------------|---------------|---------------|--------------------------------|-----------------|-------------------------|
| kubecon-aiot              | Ready  | control-plane | 6d10h | v1.21.4+k3s1 | 10.128.15.203 | 35.202.84.160 | Ubuntu 20.04.3 LTS             | 5.11.0-1018-gcp | containerd://1.4.9-k3s1 |
| follower-node-corral-tpu3 | Ready  | <none>        | 19m   | v1.21.4+k3s1 | 10.0.0.229    | <none>        | Mendel GNU/Linux 5 (Eagle)     | 4.14.98-imx     | containerd://1.4.9-k3s1 |
| follower-node-raspi       | Ready  | <none>        | 31m   | v1.21.4+k3s1 | 10.0.0.130    | <none>        | Raspbian GNU/Linux 10 (buster) | 5.10.60-v7l+    | containerd://1.4.9-k3s1 |
| follower-node-corral-tpu2 | Ready  | <none>        | 23m   | v1.21.4+k3s1 | 10.0.0.210    | <none>        | Mendel GNU/Linux 5 (Eagle)     | 4.14.98-imx     | containerd://1.4.9-k3s1 |
| follower-node-corral-tpu1 | Ready  | <none>        | 24m   | v1.21.4+k3s1 | 10.0.0.15     | <none>        | Mendel GNU/Linux 5 (Eagle)     | 4.14.98-imx     | containerd://1.4.9-k3s1 |

## Strimzi Kafka

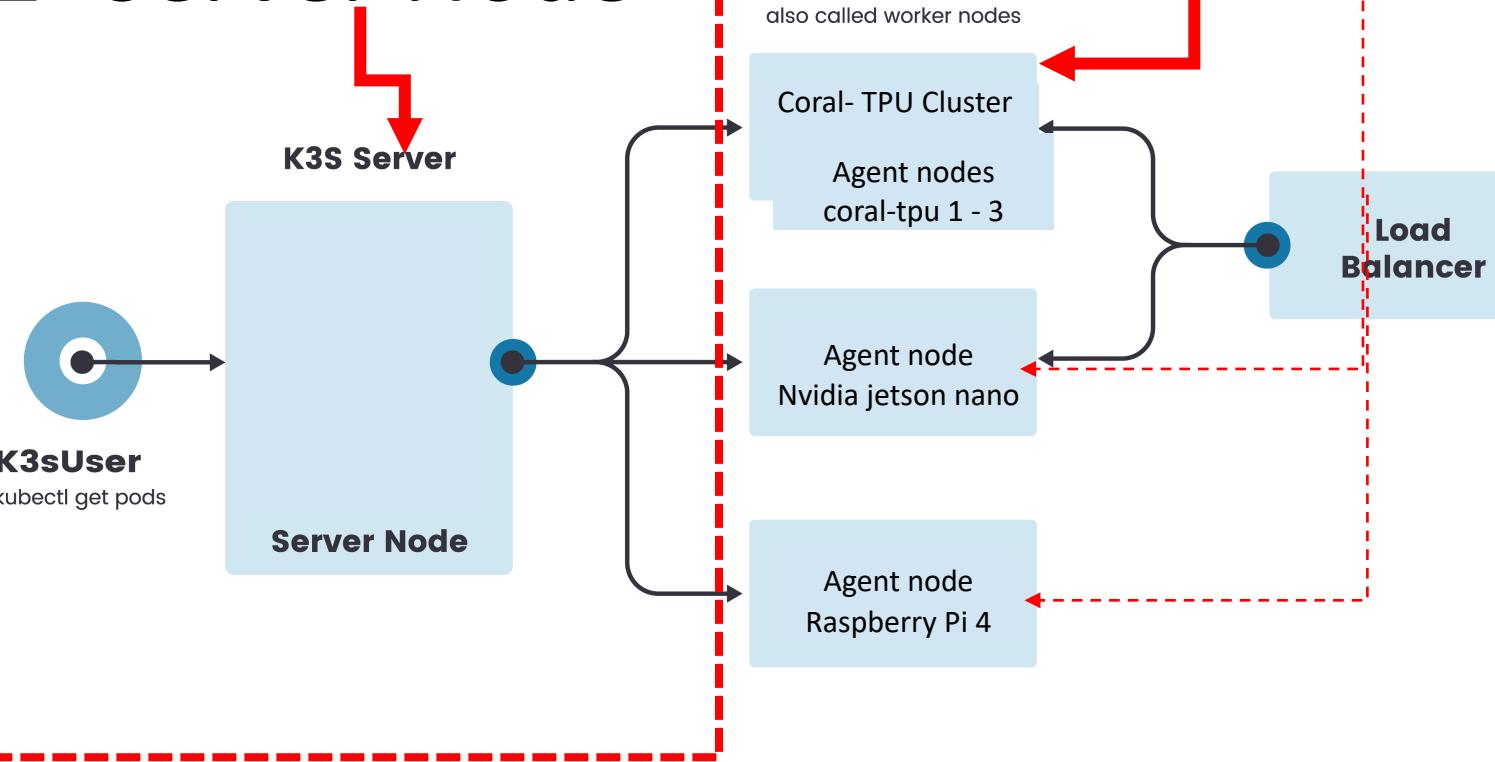
| NAME                                       | READY | STATUS  | RESTARTS | AGE   | IP         | NODE         | NOMINATED NODE | READINESS GATES |
|--|-------|---------|----------|-------|------------|--------------|----------------|-----------------|
| strimzi-cluster-operator-7695cb5f7f-5ptqk  | 1/1   | Running | 0        | 6m46s | 10.42.0.44 | kubecon-aiot | <none>         | <none>          |
| my-cluster-zookeeper-0                     | 1/1   | Running | 0        | 2m16s | 10.42.0.45 | kubecon-aiot | <none>         | <none>          |
| my-cluster-kafka-0                         | 1/1   | Running | 0        | 101s  | 10.42.0.46 | kubecon-aiot | <none>         | <none>          |
| my-cluster-entity-operator-876668cdf-hczpq | 3/3   | Running | 0        | 66s   | 10.42.0.47 | kubecon-aiot | <none>         | <none>          |

## Argo ML Ops Deployments

| NAME                                       | READY | STATUS  | RESTARTS | AGE   | IP         | NODE                      | NOMINATED NODE | READINESS GATES |
|--|-------|---------|----------|-------|------------|---------------------------|----------------|-----------------|
| strimzi-cluster-operator-7695cb5f7f-5ptqk  | 1/1   | Running | 0        | 11m   | 10.42.0.44 | kubecon-aiot              | <none>         | <none>          |
| my-cluster-zookeeper-0                     | 1/1   | Running | 0        | 6m36s | 10.42.0.45 | kubecon-aiot              | <none>         | <none>          |
| my-cluster-kafka-0                         | 1/1   | Running | 0        | 6m1s  | 10.42.0.46 | kubecon-aiot              | <none>         | <none>          |
| my-cluster-entity-operator-876668cdf-hczpq | 3/3   | Running | 0        | 5m26s | 10.42.0.47 | kubecon-aiot              | <none>         | <none>          |
| coral-python-deployment-77645bbb66-48kjq   | 2/2   | Running | 0        | 11s   | 10.42.7.7  | follower-node-corral-tpu1 | <none>         | <none>          |
| coral-python-deployment-77645bbb66-xz6qr   | 2/2   | Running | 0        | 11s   | 10.42.9.3  | follower-node-corral-tpu3 | <none>         | <none>          |
| coral-python-deployment-77645bbb66-nzch7   | 2/2   | Running | 0        | 11s   | 10.42.8.3  | follower-node-corral-tpu2 | <none>         | <none>          |

# K3S Architecture

## Step 1 - Server Node

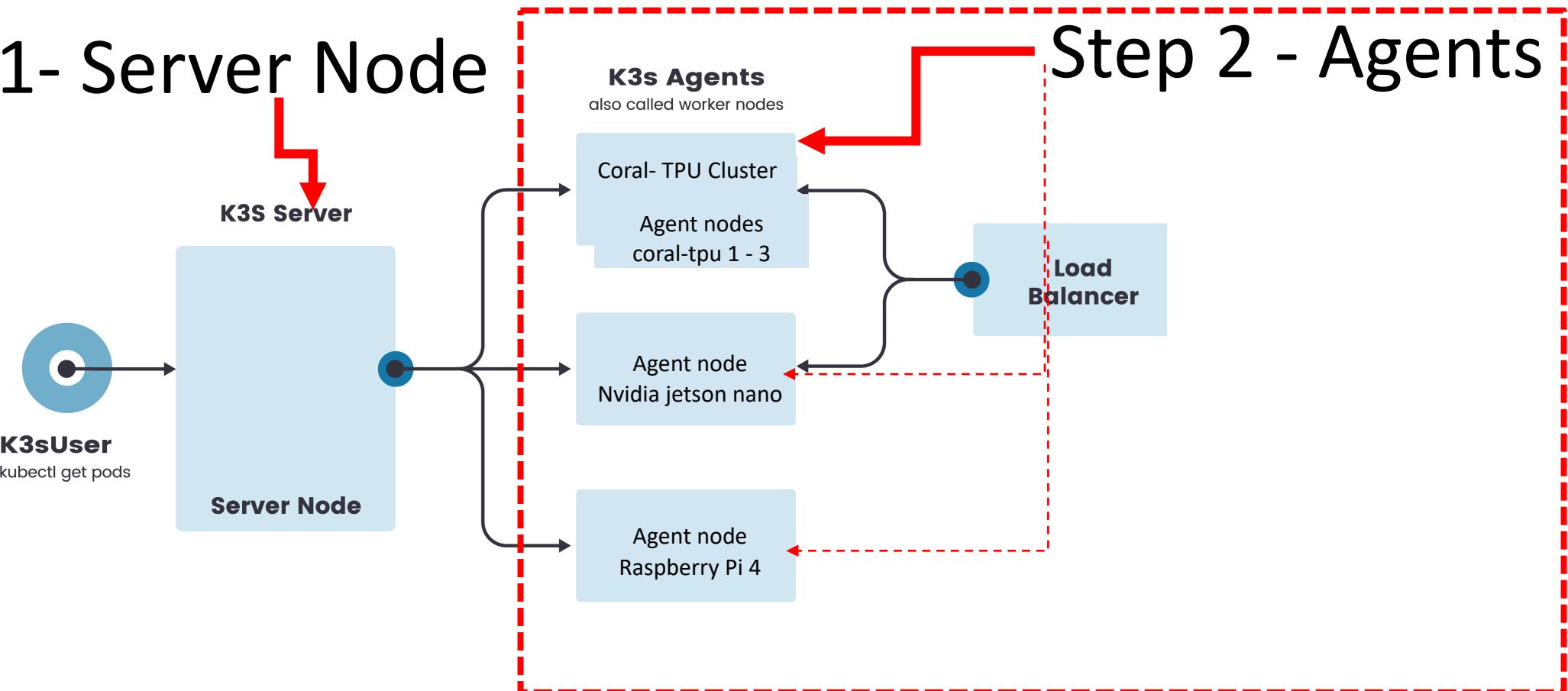


## Step 2 - Agents

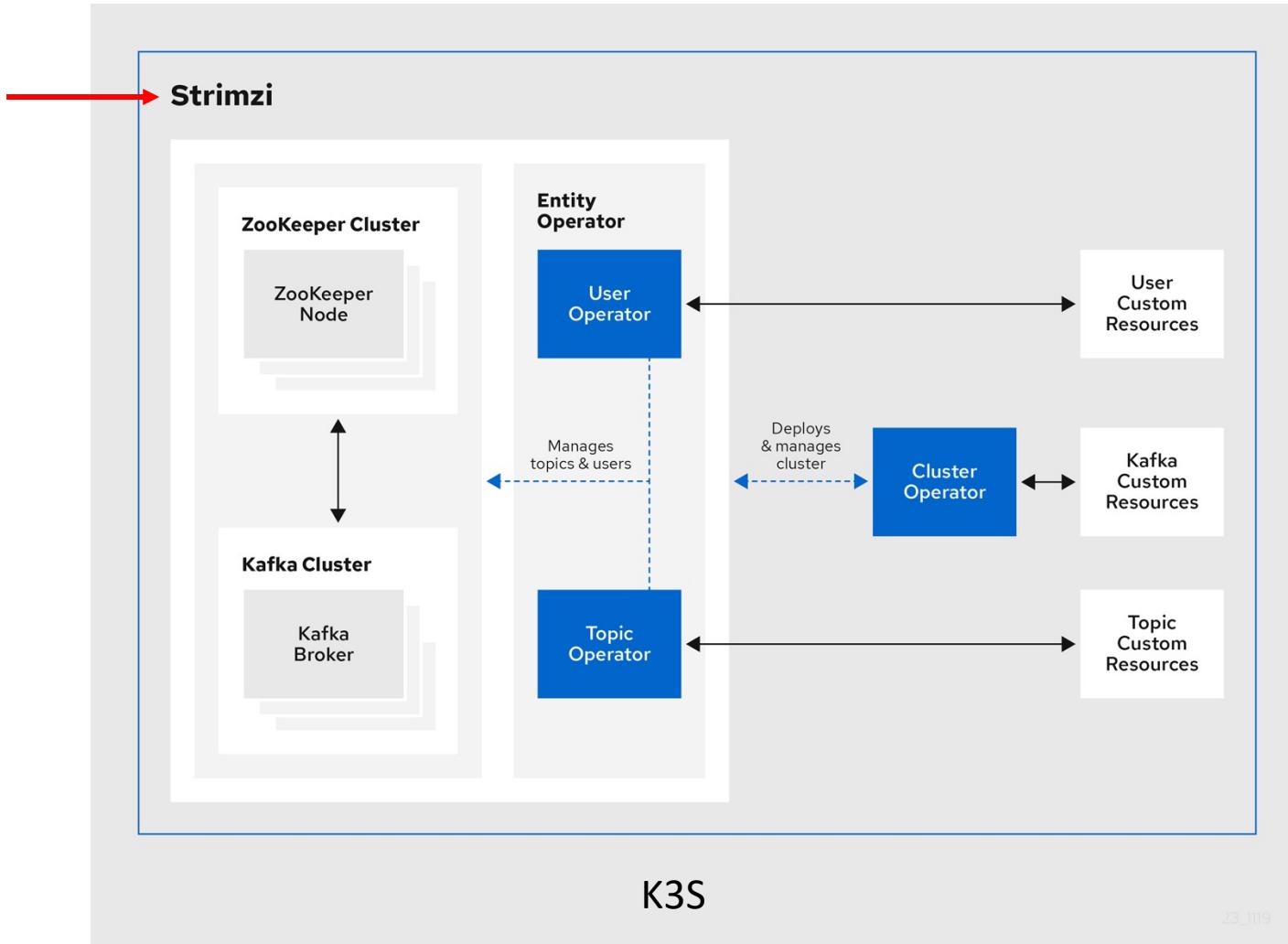
# K3S Architecture

## Step 1- Server Node

## Step 2 - Agents



# Strimzi Kafka Architecture



# Argo Workflows Architecture

