



# Using Kubernetes to Speed Development and Deployment of Edge Computer Vision Applications

Rakshit Agrawal

VP, Research & Development

Camio

# Delivering Best AI at the Edge



- From Research to Real World—Faster
- Why Containerization?
- Computer Vision in the Real World
- Confidence in Iterations
- Adaptability, Speed, and Reliability

# From Research to Real World—Faster

camio

# From Notebooks to Chipsets

```
+ Code + Text
[ ] selected_frame_crops.shape
(324, 216, 384, 3)

[ ] # Specific selection
# selected_frame_crops = get_crops_from_frames(video_frames, selections=[0,1,2,3])
# selected_frame_crops.shape

[ ] proc = [preprocess_image(i, selected_model.metadata) for i in selected_frame_crops]

[ ] np.squeeze(np.stack(proc)).shape
(324, 224, 224, 3)

[ ] mpreds = model.predict(np.squeeze(np.stack(proc)))
np.argmax(mpreds, axis=0)

array([ 0,  1, 89])

plt.hist(mpreds)
array([[120.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0., 204.],
       [207.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0., 117.],
       [321.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  3.]]),
array([0., 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1. ],
      dtype=float32),
<a list of 3 Lists of Patches objects>
300
250
200
```

ARM based edge devices



Small form appliances



Multi-socket rack servers



Data centers



# Vision AI with Containers

```
+ Code + Text
[ ] selected_frame_crops.shape
(324, 216, 384, 3)

[ ] # Specific selection
# selected_frame_crops = get_crops_from_frames(video_frames, selections=[0,1,2,3])
# selected_frame_crops.shape

[ ] proc = [preprocess_image(i, selected_model.metadata) for i in selected_frame_crops]

[ ] np.squeeze(np.stack(proc)).shape
(324, 224, 224, 3)

[ ] mpreds = model.predict(np.squeeze(np.stack(proc)))
np.argmax(mpreds, axis=0)
array([ 0,  1, 89])

[ ] plt.hist(mpreds)
array([[120.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0., 204.],
       [207.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0., 117.],
       [321.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  3.]]),
array([0., 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1. ],
      dtype=float32),
<a list of 3 Lists of Patches objects>
300
250
200
```

Low power, low memory and compute



Consumer-grade CPU/ GPU



Xeon + Nvidia RTX



Xeon + Tesla T4 grid

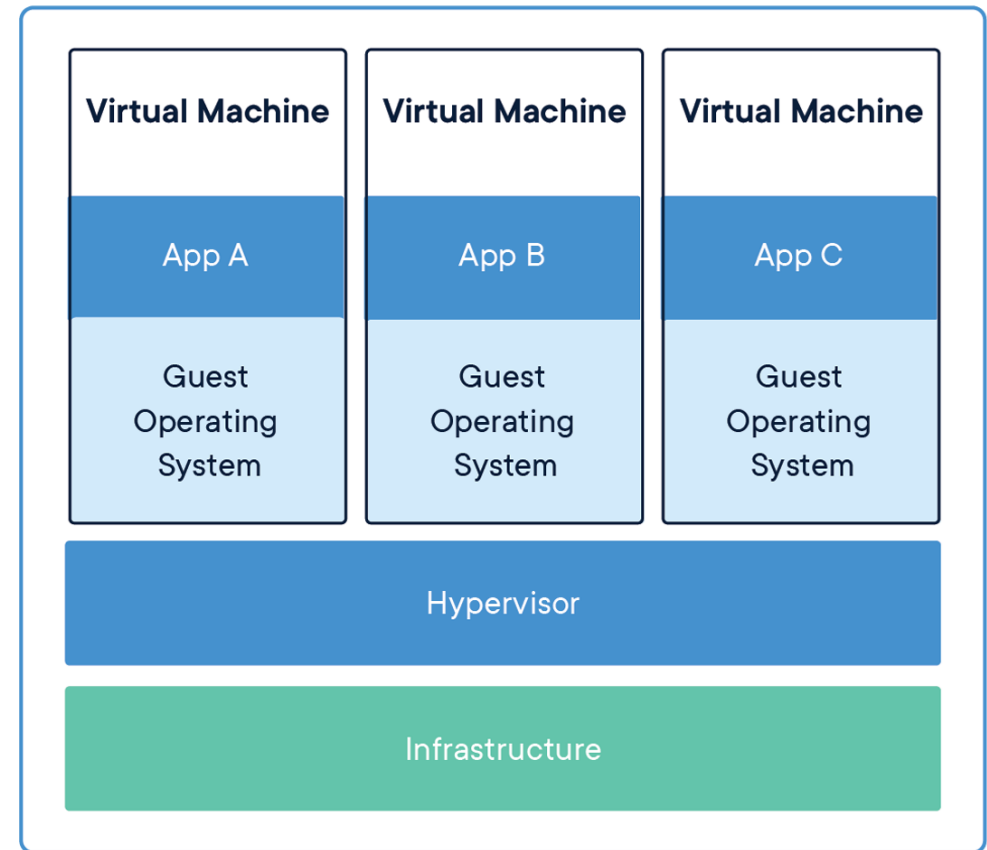
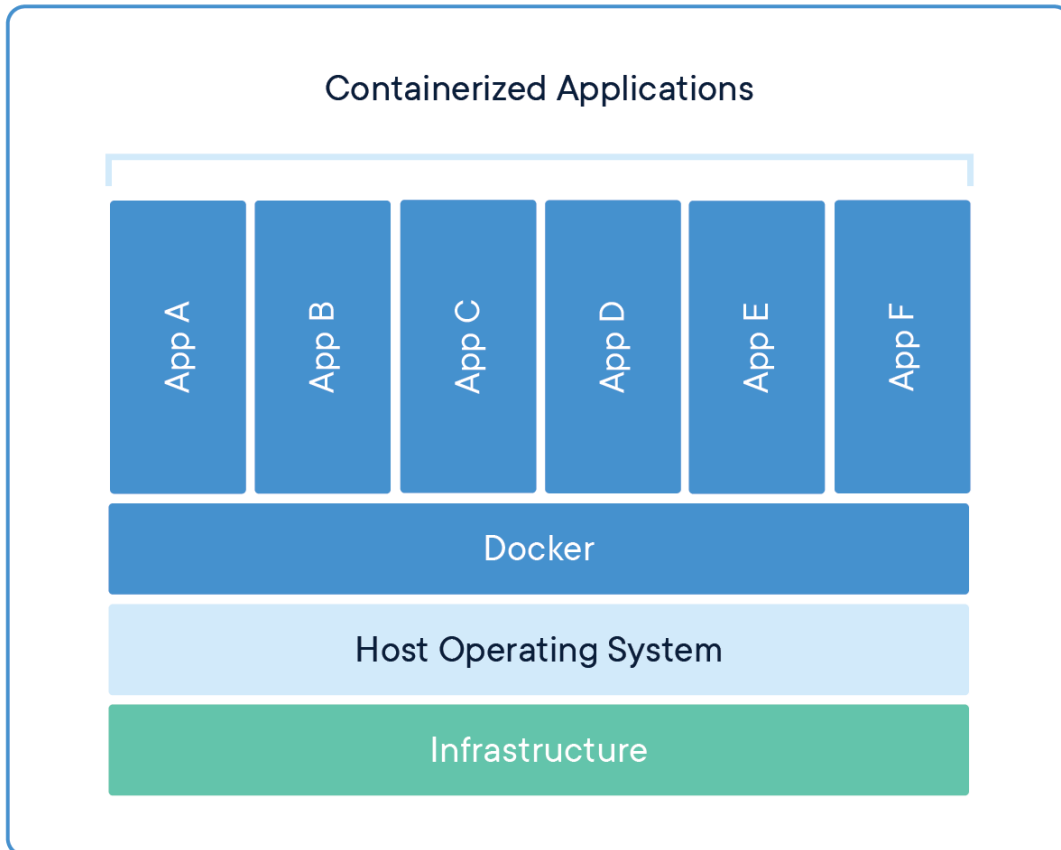


# Why Containerization?

# Fast, Flexible, Portable

- A **container** is a packaged unit of code, dependencies, and environment.
- Once built, the image can be shipped on any supporting runtime.
- A runtime or orchestrator controls deployment and lifecycle of containers.
- **Kubernetes** is a portable, extensible open-source ecosystem for orchestrating containerized workloads and services.

# Understanding Abstraction with Containers



Source: <https://www.docker.com/resources/what-container>



# Build, Configure, Deploy—Anywhere

```
FORM python:3.8
RUN apt-get update && \
    apt-get install -y sudo \
    build-essential curl \
    libcurl4-openssl-dev \
    libssl-dev wget \
    python3-pip \
    git && \
    pip3 install --upgrade pip

COPY requirements.txt .
...
```

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: dep-wq-1
  labels:
    app: dep-wq-1
spec:
  template:
    metadata:
      labels:
        app:
    ...
```

```
kubectl config get-contexts
kubectl config use-context ...

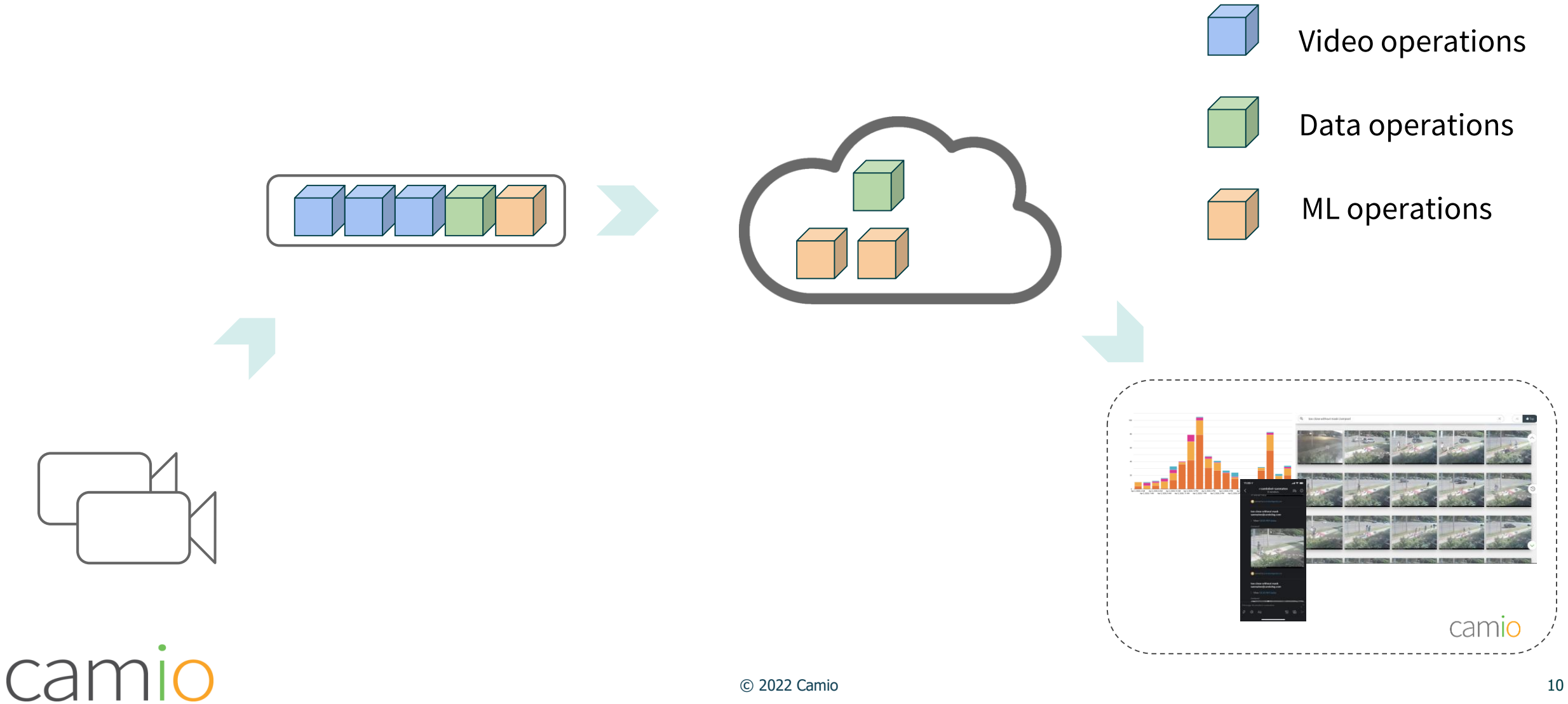
kubectl get pods

kubectl get services

kubectl apply -f pods.yaml

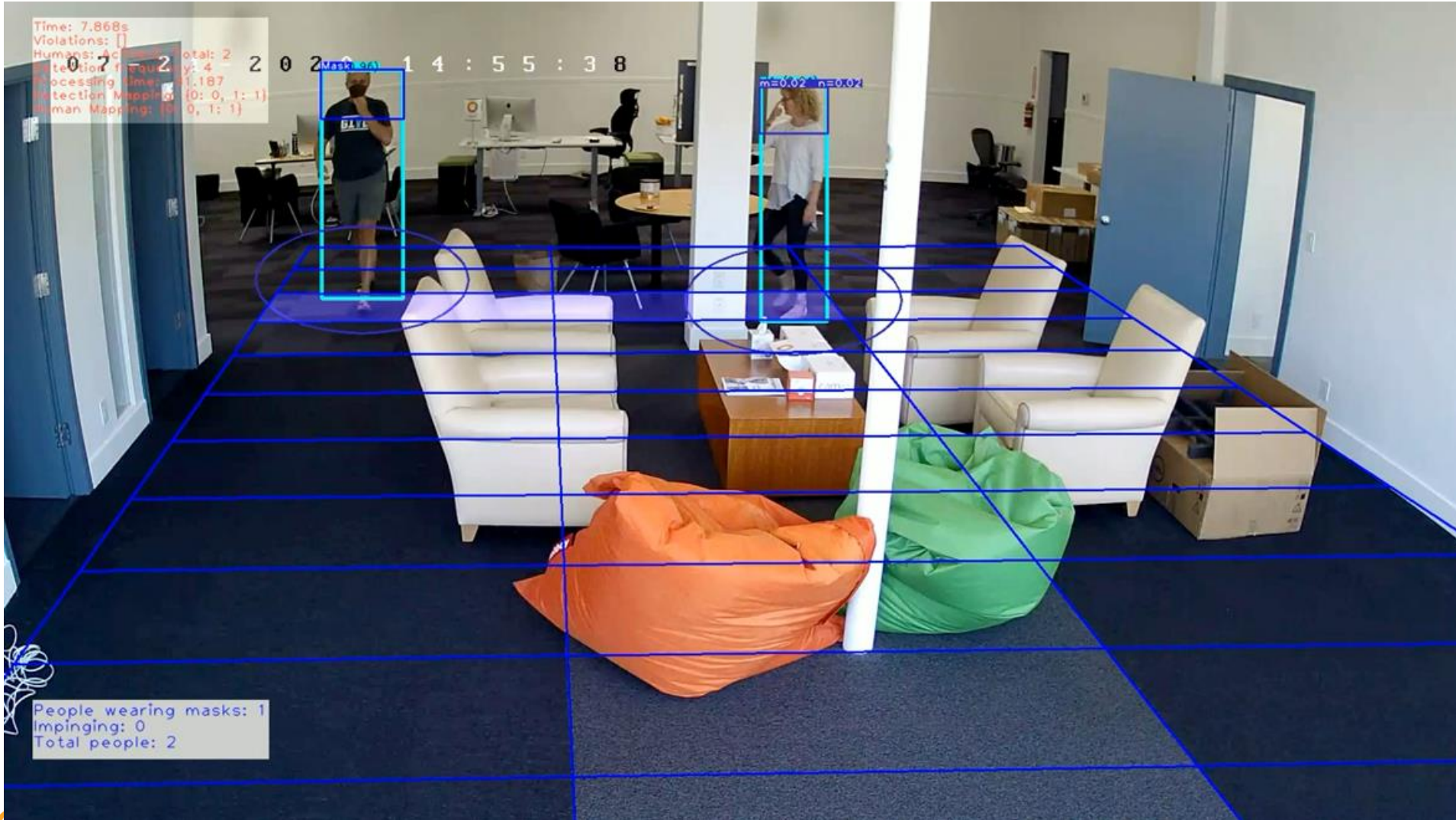
kubectl apply -f services.yaml
...
```

# Hybrid Video AI Pipeline

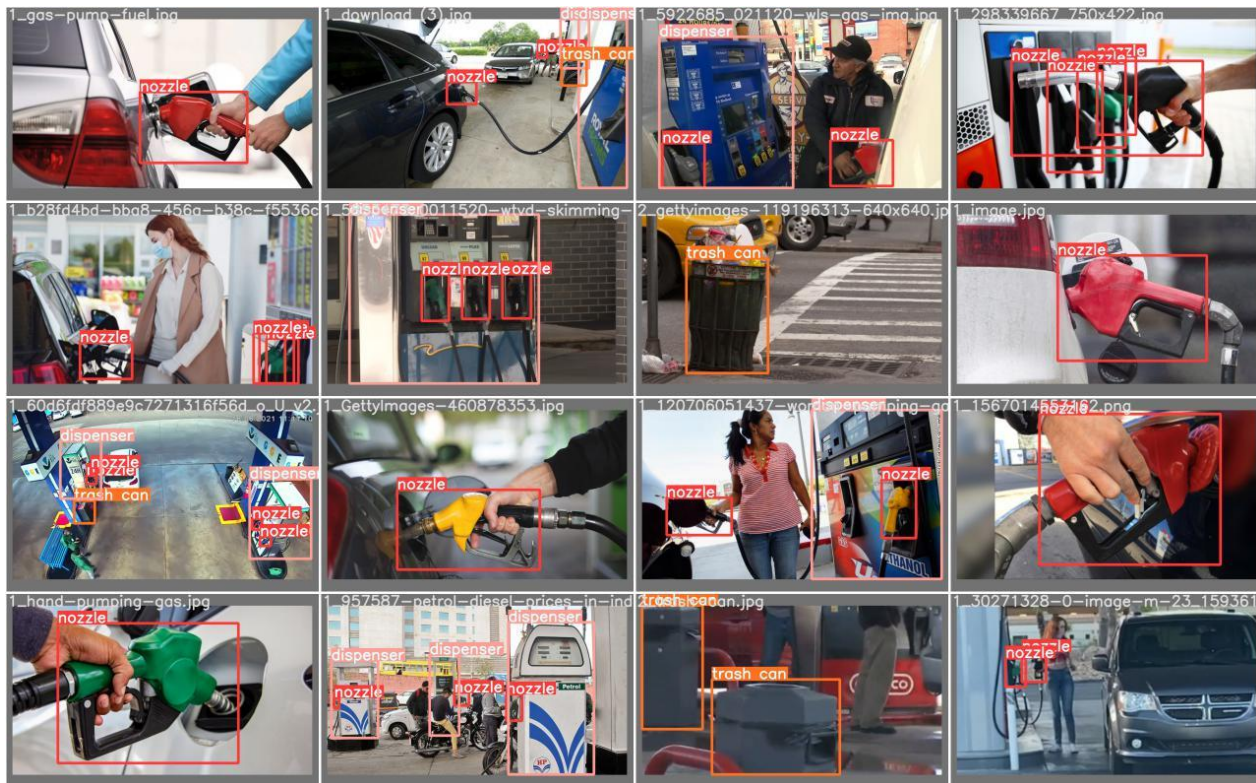


# Computer Vision in the Real World

# Video Vision AI At Work

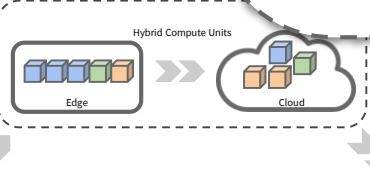
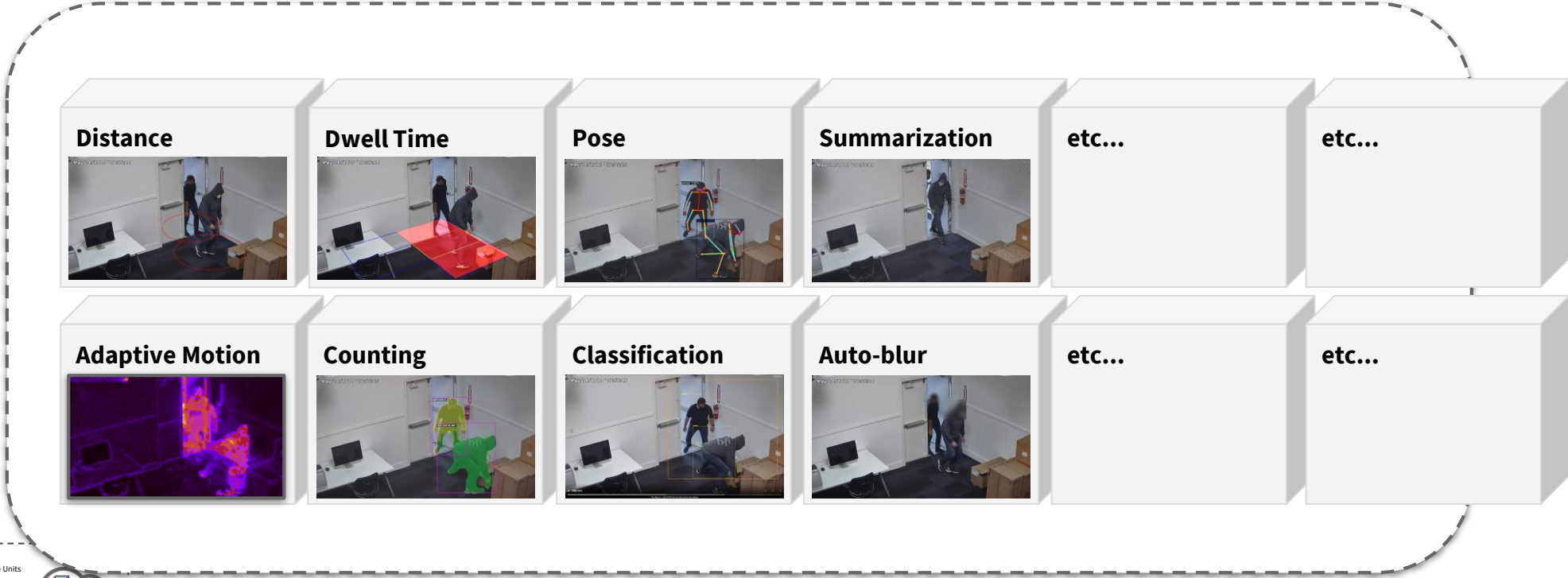


# Custom AI Deployments in One Week



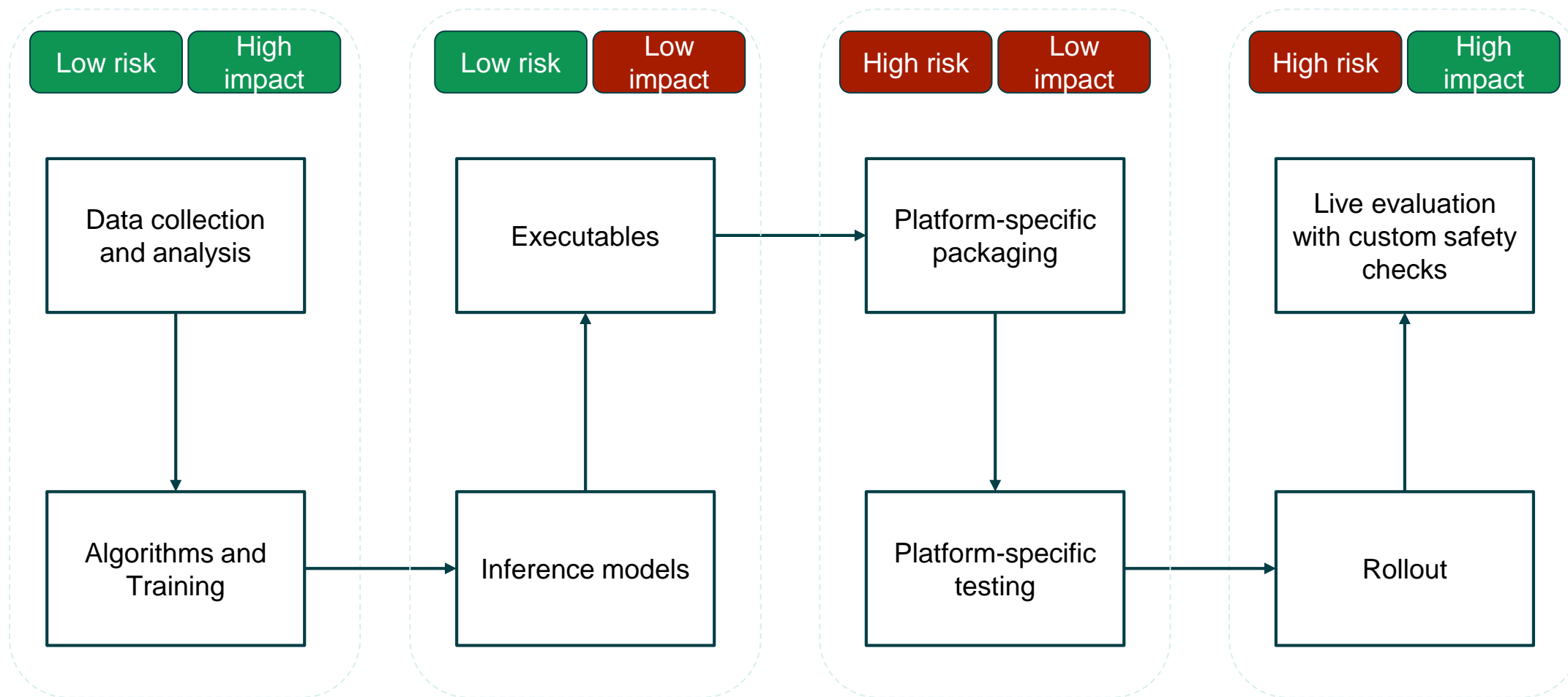
Courtesy: Gilbarco Veeder-Root

# Future-Proof systems



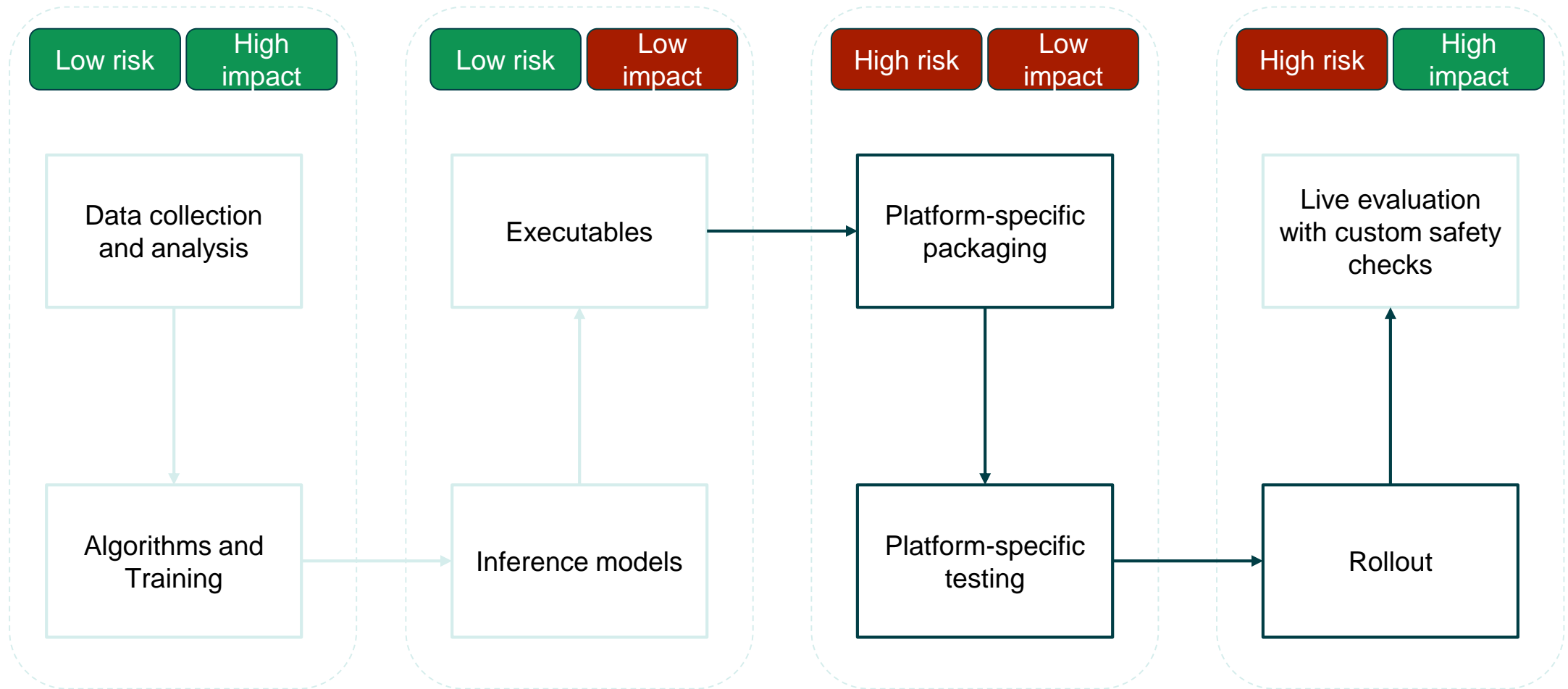
# Confidence in Iterations

# Traditional Deployment Process

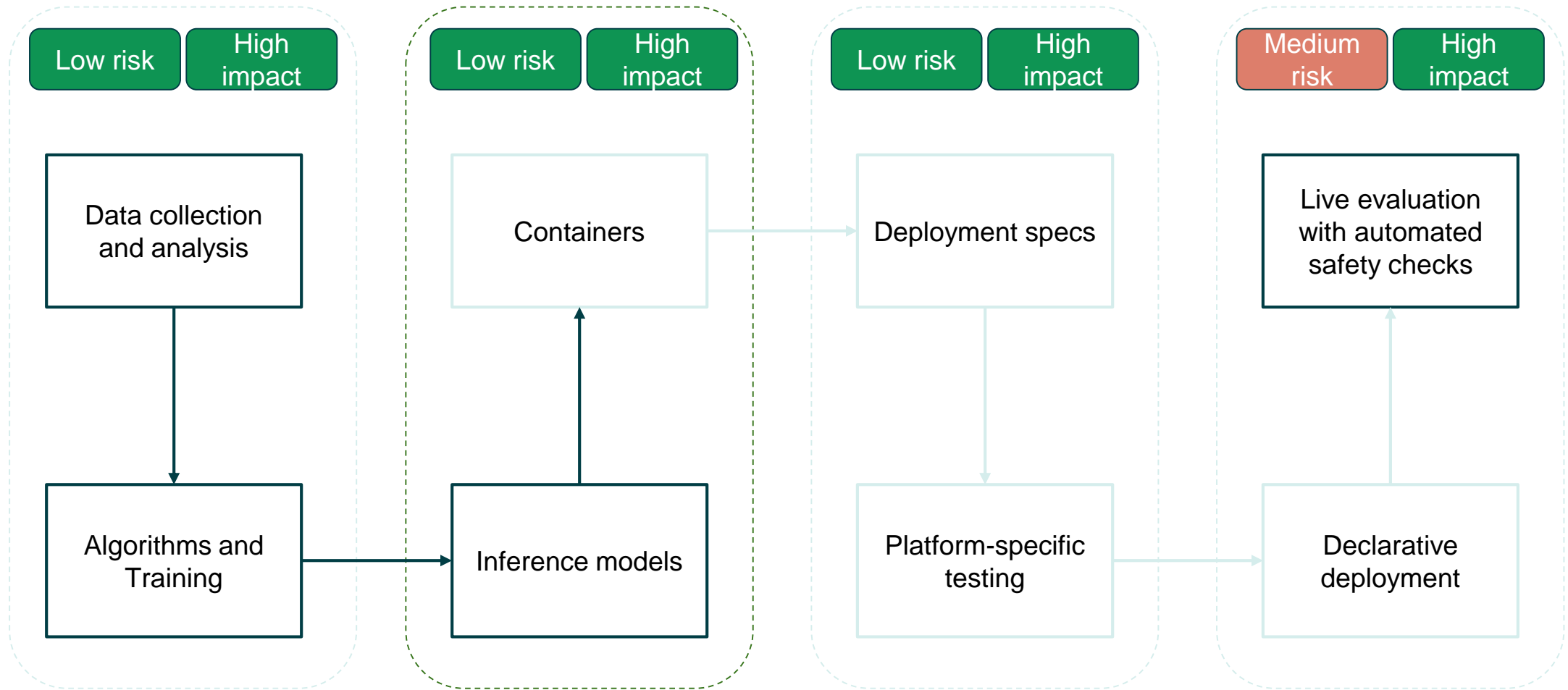




# Traditional Deployment Risk and Reward



# Containerized Deployment Process



# Key Takeaways

# Modular Infrastructure—Faster, More Flexible, More Reliable

- **Adaptability**

Containerization speeds ongoing refinements required by real-world vision AI production applications

- **Speed**

Agile pipeline operations are critical when moving from research notebooks to chipsets

- **Reliability**

- With containerization, the painful process of development, packaging and deployment becomes predictable and consistent



## Learn more about

Camio

[camio.com](https://camio.com)

Kubernetes

[kubernetes.io](https://kubernetes.io)

## Please contact for any questions or discussions

Rakshit Agrawal

Vice President, Research & Development

Camio

[rakshit@camio.com](mailto:rakshit@camio.com)