# Introduction to ML Productionization

Rachel Hu
Applied Scientist @ Amazon MLU

# Schedule

⚙ [Lecture] Intro to ML Lifecycle – 15 min

⚙ [Lab] Build an end-to-end ML system – 30 min
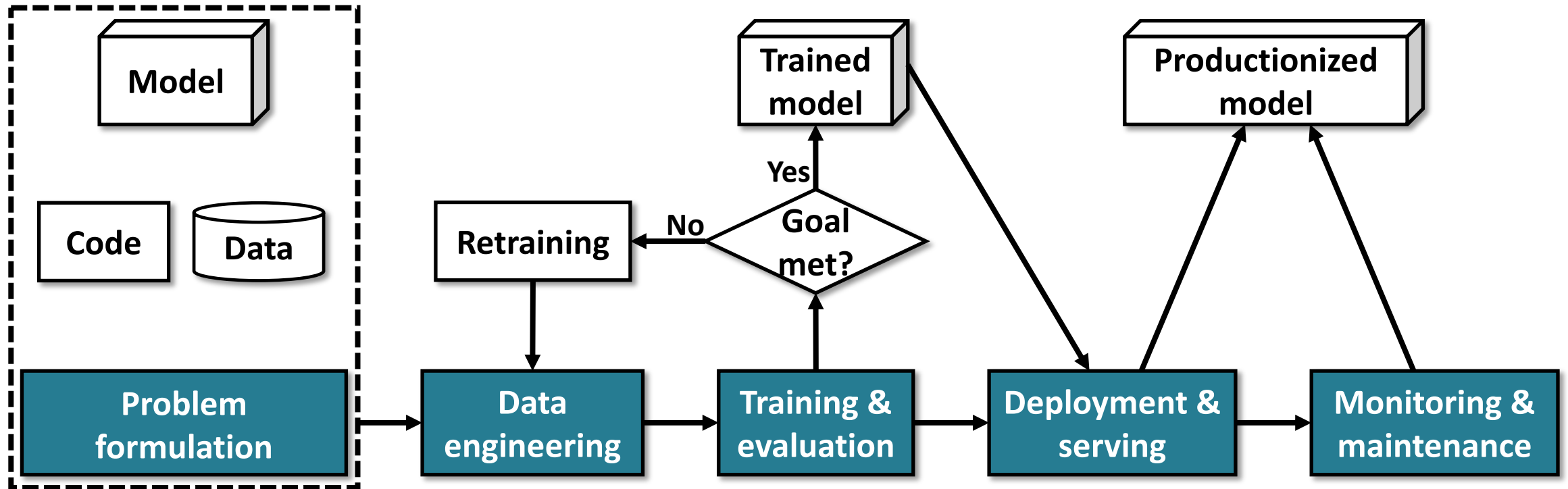  » https://github.com/goldmermaid/MLU-MLOps-Lab

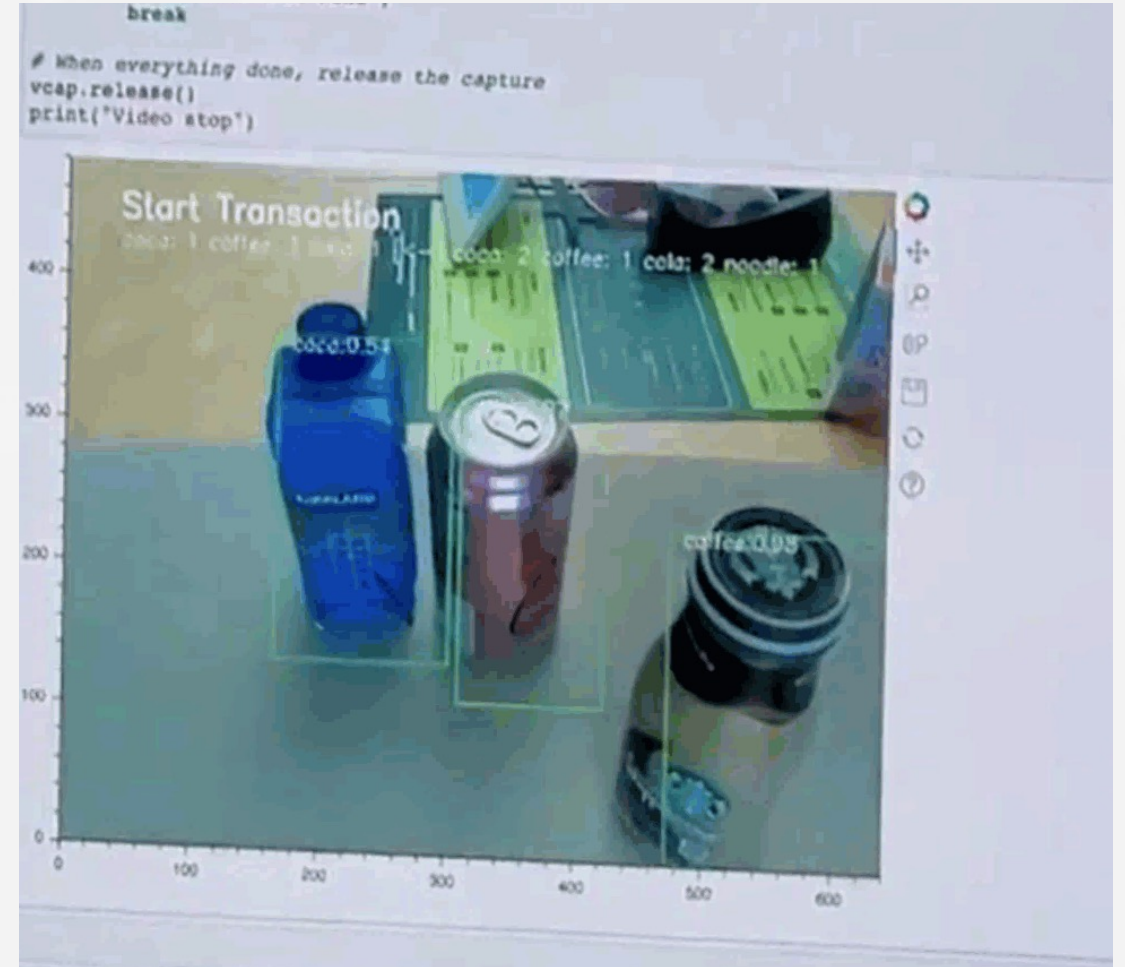⚙ Q&A – 15 min

MACHINE LEARNING UNIVERSITY

# The ML lifecycle

# A Simple ML System

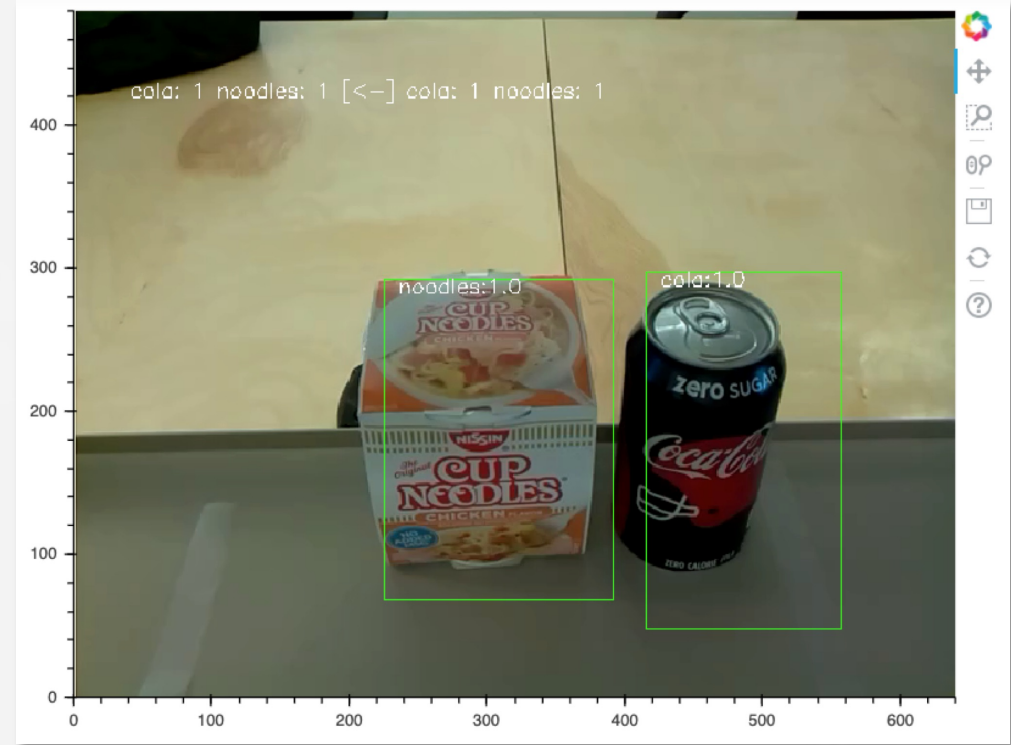Suppose that we are building a visual search system:

# A Simple ML System – Problem Formulation

1. What are the inputs and outputs?
2. How to measure this ML system's performance?
3. What are the hardware/software requirements?
   a. Batch inference, online inference or both?
   b. On-premise, cloud based or both?
   c. Which one is more cost efficient?

| Problem formulation | → | Data engineering | → | Training & evaluation | → | Deployment & serving | → | Monitoring & maintenance |
|---|---|---|---|---|---|---|---|---|

MACHINE LEARNING UNIVERSITY

# A Simple ML System - Data Engineering

⚙ Data Ingestion

  » Web sourcing

  » Simulating

  » Manually collecting

⚙ Data Processing

  » Cleaning, labeling, etc.

  » Feature engineering

  » Data augmentation



| Problem formulation | → | Data engineering | → | Training & evaluation | → | Deployment & serving | → | Monitoring & maintenance |
|---|---|---|---|---|---|---|---|---|

# A Simple ML System – Data Engineering

## Data Ingestion, Processing, and Transforming

```python
1  import boto3
2  from IPython.display import clear_output, Image, display, HTML
3  import numpy as np
4  import cv2
5  import base64
6  from bokeh.plotting import figure
7  from bokeh.io import output_notebook, show, push_notebook
8  import time
9  import json
10 output_notebook()
```

```python
1  STREAM_NAME = "pi4-001"
2  kvs = boto3.client("kinesisvideo")
3  # Grab the endpoint from GetDataEndpoint
4  endpoint = kvs.get_data_endpoint(
5      APIName="GET_HLS_STREAMING_SESSION_URL",
6      StreamName=STREAM_NAME
7  )['DataEndpoint']
8  print(endpoint)
```

• • •   • • •

```python
1  while(True):
2      ## Test frame by frame
3      ## ...
4
5      ret, frame = vcap.read()
6      if frame is not None:
7          start = time.time()
8          frame = cv2.flip(frame, -1)
9
10         # Use the trained YOLO model
11         class_IDs, scores, bounding_boxes = objectDetection.detect_image_yolo(frame)
12         frame, hand_cnt, no_hand_cnt, start_trans,
13         in_trans, curr_item_cnt, msg, msg2 = detection_result_process(
14             frame, objectDetection.classes, class_IDs, scores,
15             bounding_boxes, hand_cnt, no_hand_cnt, start_trans,
16             in_trans, curr_item_cnt, max_item_cnt, pre_msg, pre_msg2)
17
18         # Display the resulting frame
19         player(frame)
20     else:
21         print("Frame is None")
22         break
23
24 # When everything done, release the capture
25 vcap.release()
26 print("Video stop")
```

MACHINE LEARNING UNIVERSITY

# A Simple ML System – Modeling

[training_object_detector.ipynb](#)

**Training an Object Detector**

```
1  epochs = 300
2  net = train_model(train_dataset, epochs=epochs)
3  save_file = 'object_detector_epoch{}_{}.params'.format(
4      epochs, datetime.now().strftime("%m_%d_%Y_%H_%M_%S"))
5  net.save_parameters(save_file)
6  print('Saved model to disk: ' + save_file)
```

```
1   def train_model(train_dataset, epochs=50):
2       ctx = mx.gpu(0)
3       net = gcv.model_zoo.get_model('ssd_512_resnet50_v1_custom',
4                                     classes=train_dataset.classes,
5                                     transfer='coco')
6       net.collect_params().reset_ctx(ctx)
7       width, height = 512, 512   # suppose we use 512 as base training size
8       gcv.utils.random.seed(233)
9
10      batch_size = 16   # 32 for p3.2xlarge, 16 for p2.2xlarge
11      num_workers = 4
12      with autograd.train_mode():
13          _, _, anchors = net(mx.nd.zeros((1, 3, height, width), ctx))
14      anchors = anchors.as_in_context(mx.cpu())
15      train_transform = SSDDefaultTrainTransform(width, height, anchors)
16      batchify_fn = Tuple(Stack(), Stack(), Stack())
17      train_loader = mx.gluon.data.DataLoader(
18          train_dataset.transform(train_transform),
19          batch_size,
20          shuffle=True,
21          batchify_fn=batchify_fn,
22          last_batch='rollover',
23          num_workers=num_workers)
```
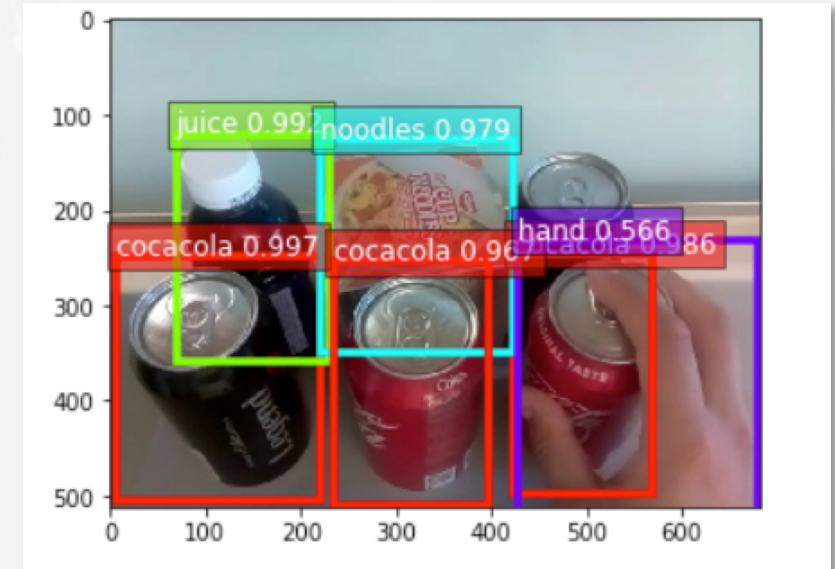
• • •   • • •

# A Simple ML System - Testing

[test_object_detector.ipynb](test_object_detector.ipynb)

## Validating a Model

```python
1   def validate(net, test_dataset, ctx):
2       if isinstance(ctx, mx.Context):
3           ctx = [ctx]
4       size = len(test_dataset)
5       metric = gcv.utils.metrics.voc_detection.VOC07MApMetric(
6           iou_thresh=0.5, class_names=test_dataset.classes)
7       net.collect_params().reset_ctx(ctx)
8       metric.reset()
9       width, height = 512, 512
10      batch_size = 4
11      batchify_fn = Tuple(Stack(), Pad(pad_val=-1))
12      val_loader = mx.gluon.data.DataLoader(
13          test_dataset.transform(SSDDefaultValTransform(width, height)),
14          batchify_fn=batchify_fn, batch_size=batch_size, shuffle=False,
15          last_batch='rollover', num_workers=0)
```
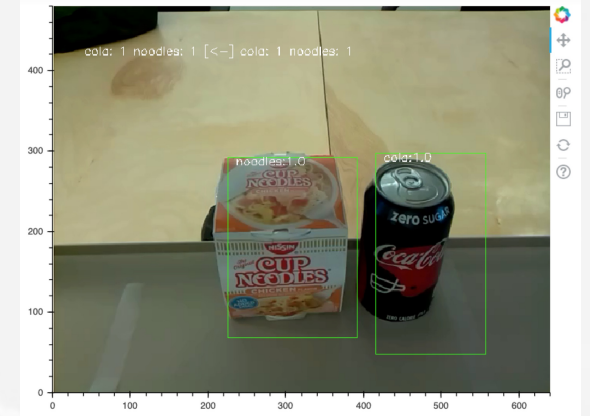
... ...

# Deployment & Serving



```python
while(True):
    ## Test frame by frame
    ## ...

    ret, frame = vcap.read()
    if frame is not None:
        start = time.time()
        frame = cv2.flip(frame, -1)

        # Use the trained YOLO model
        class_IDs, scores, bounding_boxes = objectDetection.detect_image_yolo(frame)
        frame, hand_cnt, no_hand_cnt, start_trans,
        in_trans, curr_item_cnt, msg, msg2 = detection_result_process(
            frame, objectDetection.classes, class_IDs, scores,
            bounding_boxes, hand_cnt, no_hand_cnt, start_trans,
            in_trans, curr_item_cnt, max_item_cnt, pre_msg, pre_msg2)

        # Display the resulting frame
        player(frame)
    else:
        print("Frame is None")
        break

# When everything done, release the capture
vcap.release()
print("Video stop")
```
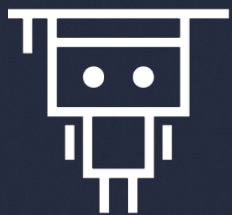
VideoStream.ipydb

```python
class ObjectDetection():
    def __init__(self):
        self.classes = ['cocacola', 'juice', 'noodles', 'hand']  # , 'cocacola-zero'
        self.net = model_zoo.get_model('ssd_512_resnet50_v1_custom',
                                       classes=self.classes, pretrained_base=False)
#        self.net = model_zoo.get_model('yolo3_darknet53_custom',
#                                       classes=self.classes, pretrained_base=False)
        param_files = ([x for x in os.listdir('.') if x.endswith('.params')])
        selected = param_files[0]
        self.net.load_parameters(selected)
        self.ctx = mx.gpu(0)
        self.net.collect_params().reset_ctx(self.ctx)
```
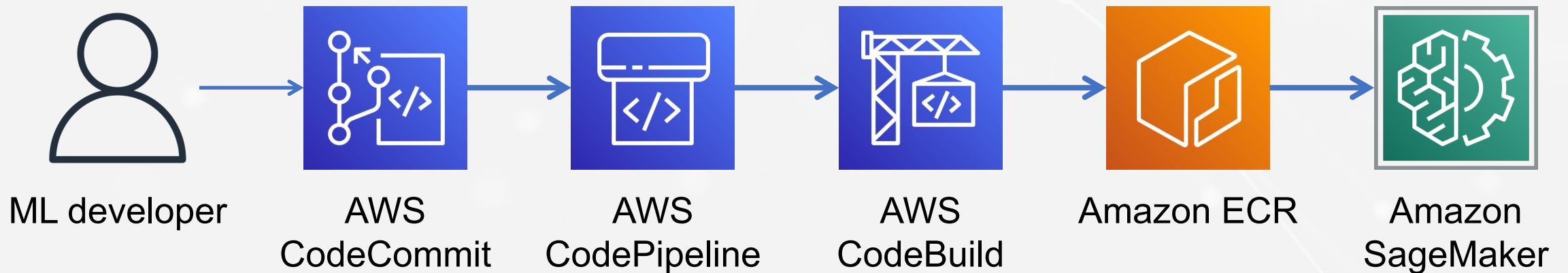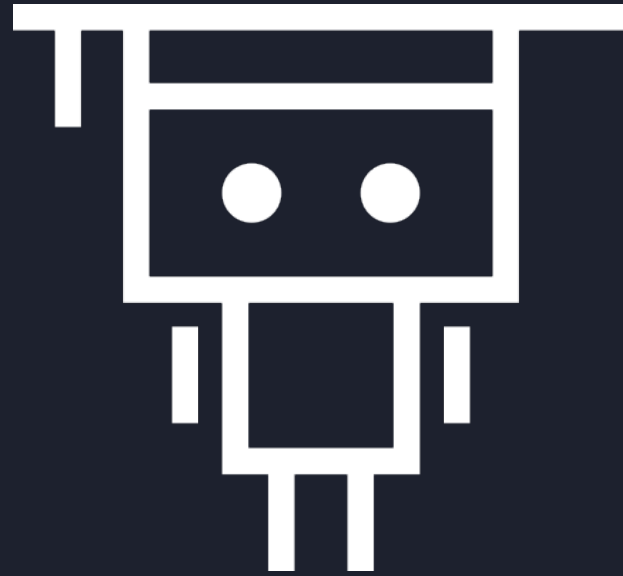
object_detection.py

MACHINE LEARNING UNIVERSITY

# Lab

# Lab Exercise

In this exercise, we will create an MLOps pipeline using CloudFormation



ML developer → AWS CodeCommit → AWS CodePipeline → AWS CodeBuild → Amazon ECR → Amazon SageMaker

Lab instructions: https://github.com/goldmermaid/MLU-MLOps-Lab

MACHINE LEARNING UNIVERSITY

- Email: rlhu@amazon.com
- LinkedIn: https://www.linkedin.com/in/rachelsonghu

# Thank you!