AWS re:Invent

AIM404

# Build, Train, and Deploy ML Models with Amazon SageMaker

David Arpin
Data Science Practice Manager
AWS, Professional Services

**Featuring:** Prasad Prabhu
Principal Architect
Intuit, Data Platform

aws re:Invent

aws

# Agenda

## Review Amazon SageMaker
- Build, train, and deploy
- Algorithms, frameworks, bring your own, and automatic model tuning

## Realtime deployment at scale
- Creating and updating endpoints
- Reduced risk deployments
- Automatic scaling

## Customer story
- ML at Intuit
- Data science workflows
- Architecture and demo

AWS re:Invent

aws

# The Amazon Machine Learning Stack

## AI SERVICES

| | Vision | | Speech | | Language | | Chatbots & Contact Centers |
|---|---|---|---|---|---|---|---|
| | AMAZON REKOGNITION IMAGE | AMAZON REKOGNITION VIDEO | AMAZON POLLY | AMAZON TRANSCRIBE | AMAZON TRANSLATE | AMAZON COMPREHEND | AMAZON LEX |

## ML SERVICES

AMAZON SAGEMAKER

## ML FRAMEWORKS & INFRASTRUCTURE

### Frameworks

TensorFlow   mxnet   PYTORCH   Chainer   HOROVOD

### Interfaces

GLUON   K Keras

### Infrastructure

AMAZON EC2 P3 Instances     AMAZON EC2 C5 Instances     FPGAs

AWS re:Invent

aws

# ML is still too complicated for everyday developers

Collect and prepare training data
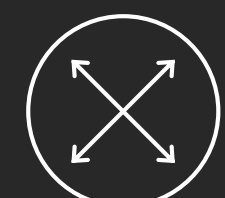
Choose and optimize your ML algorithm

Set up and manage environments for training

Train and tune model (trial and error)

Deploy model in production

Scale and manage the production environment

AWS re:Invent

aws

# Machine Learning Made Simple

AMAZON
SAGEMAKER

**One**-click

model training
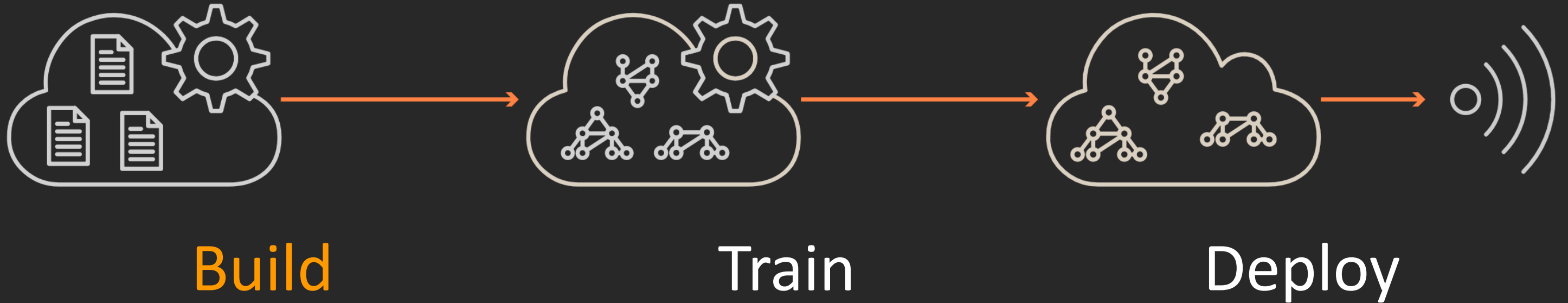& deployment

**10x**

better algorithm
performance

**Predictive** insights

to improve business
decision making

AWS
re:Invent

aws

# Amazon SageMaker simplifies Machine Learning



Build → Train → Deploy

# Amazon SageMaker modules



**Build**             Train             Deploy

- Notebook instances
- Call APIs from your device

# Amazon SageMaker modules
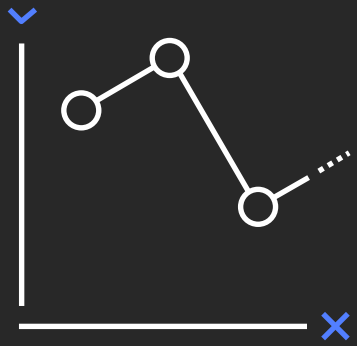
Build ⟶ Train ⟶ Deploy

- Managed
- Distributed
- High performance I/O

# Amazon SageMaker modules

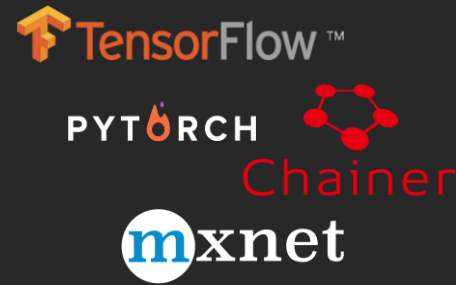Build       Train       Deploy

- Real-time endpoints
- Batch transform
- AWS Greengrass
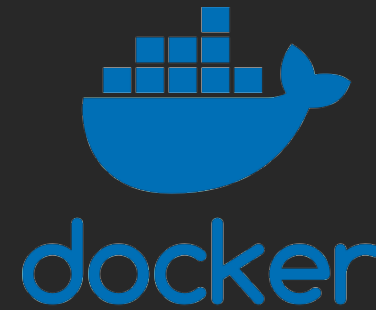- AWS DeepLens
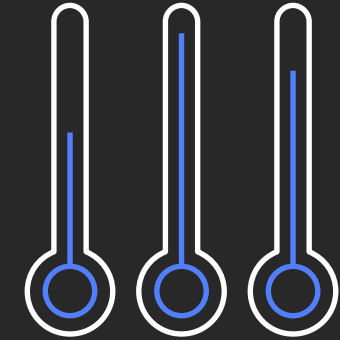
aws

# Amazon SageMaker features



SageMaker

algorithms

Frameworks

Bring

your own

Automatic

model tuning

aws

# Amazon SageMaker features



**SageMaker algorithms**

- Designed for speed and scale
- Supervised, unsupervised, computer vision, and NLP



Frameworks



Bring your own



Automatic model tuning

# Amazon SageMaker features

SageMaker
algorithms

**Frameworks**

- 20 lines of Python
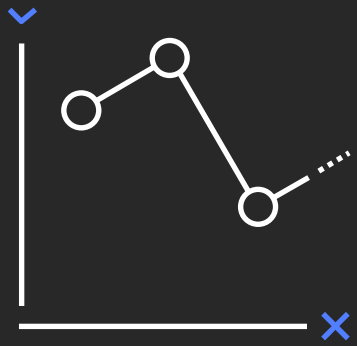- Open sourced
- Local mode for testing

Bring

your own

Automatic

model tuning

# Amazon SageMaker features



SageMaker
algorithms



Frameworks



**Bring**

**your own**

- Publish to a container registry
- R, Java, Julia, etc.



Automatic
model tuning

# Amazon SageMaker features
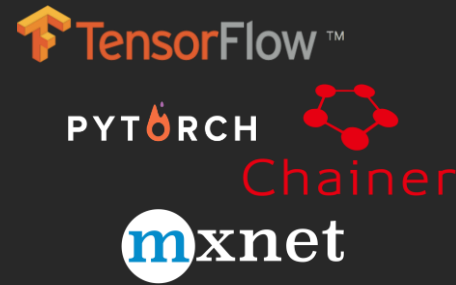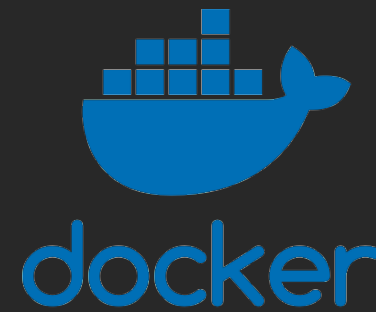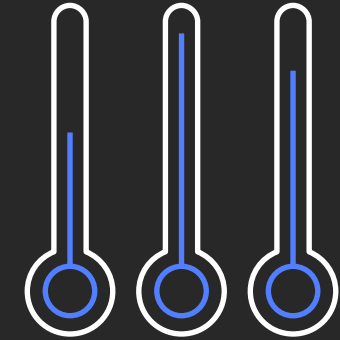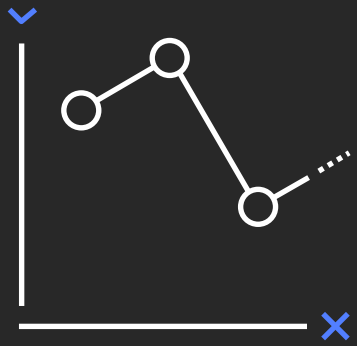
SageMaker algorithms

Frameworks

Bring your own

## Automatic model tuning

- Efficient meta-model hyperparameter tuning
- Works with algorithms, frameworks, and BYO

# Real-time deployment at scale

# Creating endpoints

Easy deployment to production REST API

Scalable, high throughput, and high reliability

# Creating endpoints

Model

```
aws sagemaker create-model
    --model-name model1
    --primary-container '{"Image": "123.dkr.ecr.amazonaws.com/algo",
                          "ModelDataUrl": "s3://bkt/model1.tar.gz"}'
    --execution-role-arn arn:aws:iam::123:role/me
```

# Creating endpoints

Model

```
aws sagemaker create-model
    --model-name model1
    --primary-container '{"Image": "123.dkr.ecr.amazonaws.com/algo",
                          "ModelDataUrl": "s3://bkt/model1.tar.gz"}'
    --execution-role-arn arn:aws:iam::123:role/me
```

Endpoint

configuration

```
aws sagemaker create-endpoint-config
    --endpoint-config-name model1-config
    --production-variants '{"InitialInstanceCount": 2,
                            "InstanceType": "ml.m4.xlarge",
                            "InitialVariantWeight": 1,
                            "ModelName": "model1",
                            "VariantName": "AllTraffic"}'
```

aws

# Creating endpoints

Model

```
aws sagemaker create-model
    --model-name model1
    --primary-container '{"Image": "123.dkr.ecr.amazonaws.com/algo",
                          "ModelDataUrl": "s3://bkt/model1.tar.gz"}'
    --execution-role-arn arn:aws:iam::123:role/me
```

Endpoint
configuration

```
aws sagemaker create-endpoint-config
    --endpoint-config-name model1-config
    --production-variants '{"InitialInstanceCount": 2,
                            "InstanceType": "ml.m4.xlarge",
                            "InitialVariantWeight": 1,
                            "ModelName": "model1",
                            "VariantName": "AllTraffic"}'
```

Endpoint

```
aws sagemaker create-endpoint
    --endpoint-name my-endpoint
    --endpoint-config-name model1-config
```

# Updating endpoints

Blue-green deployments mean no scheduled downtime

Deploy one or more models behind the same endpoint

# Updating endpoints

New model

```
aws sagemaker create-model
    --model-name model2
    --primary-container '{"Image": "123.dkr.ecr.amazonaws.com/algo",
                          "ModelDataUrl": "s3://bkt/model2.tar.gz"}'
    --execution-role-arn arn:aws:iam::123:role/me
```

# Updating endpoints

**New model**

```
aws sagemaker create-model
    --model-name model2
    --primary-container '{"Image": "123.dkr.ecr.amazonaws.com/algo",
                          "ModelDataUrl": "s3://bkt/model2.tar.gz"}'
    --execution-role-arn arn:aws:iam::123:role/me
```

**New endpoint**

**configuration**

```
aws sagemaker create-endpoint-config
    --endpoint-config-name model2-config
    --production-variants '{"InitialInstanceCount": 2,
                            "InstanceType": "ml.m4.xlarge",
                            "InitialVariantWeight": 1,
                            "ModelName": "model2",
                            "VariantName": "AllTraffic"}'
```

# Updating endpoints

New model

```
aws sagemaker create-model
    --model-name model2
    --primary-container '{"Image": "123.dkr.ecr.amazonaws.com/algo",
                          "ModelDataUrl": "s3://bkt/model2.tar.gz"}
    --execution-role-arn arn:aws:iam::123:role/me
```
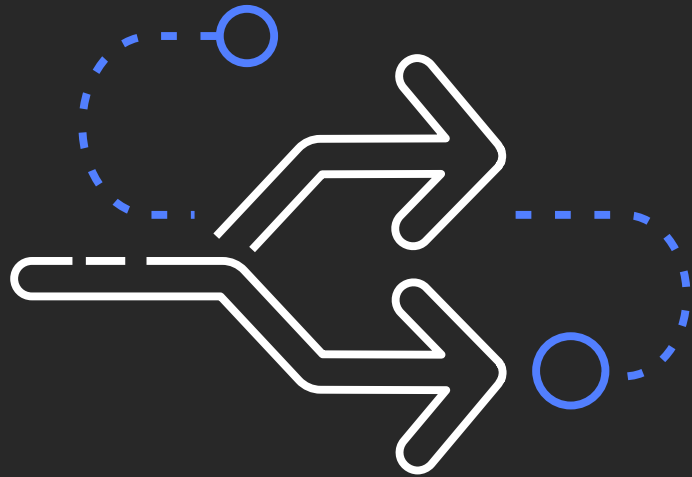
New endpoint

configuration

```
aws sagemaker create-endpoint-config
    --endpoint-config-name model2-config
    --production-variants '{"InitialInstanceCount": 2,
                            "InstanceType": "ml.m4.xlarge",
                            "InitialVariantWeight": 1,
                            "ModelName": "model2",
                            "VariantName": "AllTraffic"}'
```
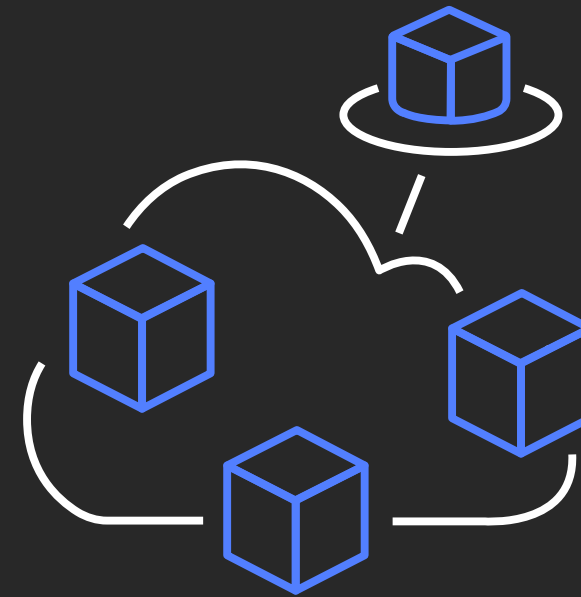
Same

endpoint

```
aws sagemaker update-endpoint
    --endpoint-name my-endpoint
    --endpoint-config-name model2-config
```

aws

# Reduced risk deployments

Incrementally retrain
models with new data

Try new models and
improved algorithms

# Reduced risk deployments

Two-model
endpoint
configuration

```
aws sagemaker create-endpoint-config
    --endpoint-config-name both-models-config
    --production-variants '[{"InitialInstanceCount": 2,
                "InstanceType": "ml.m4.xlarge",
                "InitialVariantWeight": 95,
                "ModelName": "model1",
                "VariantName": "model1-traffic"},
                {"InitialInstanceCount": 2,
                "InstanceType": "ml.m4.xlarge",
                "InitialVariantWeight": 5,
                "ModelName": "model2",
                "VariantName": "model2-traffic"}]'
```
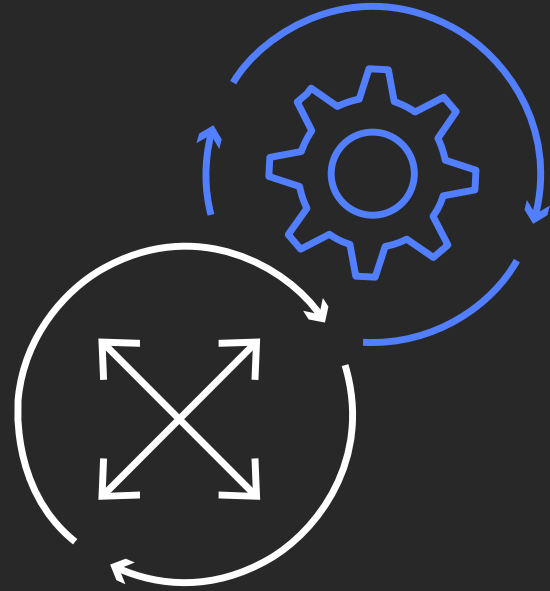
# Reduced risk deployments

Two-model
endpoint
configuration

```
aws sagemaker create-endpoint-config
    --endpoint-config-name both-models-config
    --production-variants '[{"InitialInstanceCount": 2,
                            "InstanceType": "ml.m4.xlarge",
                            "InitialVariantWeight": 95,
                            "ModelName": "model1",
                            "VariantName": "model1-traffic"},
                           {"InitialInstanceCount": 2,
                            "InstanceType": "ml.m4.xlarge",
                            "InitialVariantWeight": 5,
                            "ModelName": "model2",
                            "VariantName": "model2-traffic"}]'
```
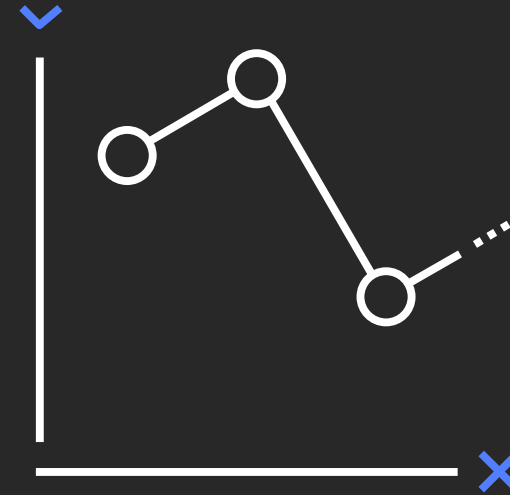
Same
endpoint

```
aws sagemaker update-endpoint
    --endpoint-name my-endpoint
    --endpoint-config-name both-models-config
```

aws

# Reduced risk deployments

**Two-model endpoint configuration**

```
aws sagemaker create-endpoint-config
    --endpoint-config-name both-models-config
    --production-variants '[{"InitialInstanceCount": 2,
                            "InstanceType": "ml.m4.xlarge",
                            "InitialVariantWeight": 95,
                            "ModelName": "model1",
                            "VariantName": "model1-traffic"},
                          {"InitialInstanceCount": 2,
                            "InstanceType": "ml.m4.xlarge",
                            "InitialVariantWeight": 5,
                            "ModelName": "model2",
                            "VariantName": "model2-traffic"}]'
```

**Same endpoint**

```
aws sagemaker update-endpoint
    --endpoint-name my-endpoint
    --endpoint-config-name both-models-config
```

**Swap**

```
aws sagemaker update-endpoint-weights-and-capacities
    --endpoint-name my-endpoint
    --desired-weights-and-capacities '{"VariantName": "model1",
                                       "DesiredWeight": 5}'
```

# Automatic scaling endpoints

SageMaker console settings:

- Min and max instances

- Target invocations per instance

- Scaling cooldowns

---

**Variant automatic scaling** Learn more

| Variant name | Instance type | Current instance count | Current weight |
|---|---|---|---|
| AllTraffic | ml.p2.xlarge | 2 | 1 |

Minimum instance count    Maximum instance count

| 2 | - | 5 |

IAM role
Amazon SageMaker uses the following service-linked role for automatic scaling. Learn more

AWSServiceRoleForApplicationAutoScaling_SageMakerEndpoint

**Built-in scaling policy** Learn more

Policy name

SageMakerEndpointInvocationScalingPolicy

Target metric                                    Target value

SageMakerVariantInvocationsPerInstance 🔗    | 800 |

Scale in cool down (seconds) - *optional*    Scale out cool down (seconds) - *optional*

| 120 |    | 60 |

AWS re:Invent

aws

# Why automatic scaling?

# Automatic scaling in action

# Scaling criteria

Algorithms have

different memory, CPU,

or GPU requirements

Automatically scale

based on endpoint

instance's Amazon

CloudWatch metrics

# Creating an automatic scaling policy

Variant

```
aws application-autoscaling register-scalable-target
    --service-namespace sagemaker
    --resource-id endpoint/my-endpoint/variant/model2
    --scalable-dimension sagemaker:variant:DesiredInstanceCount
    --min-capacity 2
    --max-capacity 5
```

aws

# Creating an automatic scaling policy

Variant

```
aws application-autoscaling register-scalable-target
    --service-namespace sagemaker
    --resource-id endpoint/my-endpoint/variant/model2
    --scalable-dimension sagemaker:variant:DesiredInstanceCount
    --min-capacity 2
    --max-capacity 5
```

Policy

```
aws application-autoscaling put-scaling-policy
    --policy-name model2-scaling
    --service-namespace sagemaker
    --resource-id endpoint/my-endpoint/variant/model2
    --scalable-dimension sagemaker:variant:DesiredInstanceCount
    --policy-type TargetTrackingScaling
    --target-tracking-scaling-policy-configuration
    '{"TargetValue": 50,
        "CustomizedMetricSpecification":
            {"MetricName": "CPUUtilization",
             "Namespace": "/aws/sagemaker/Endpoints",
             "Dimensions":
                [{"Name": "EndpointName", "Value": "my-endpoint"},
                 {"Name": "VariantName","Value": "model2"}],
             "Statistic": "Average",
             "Unit": "Percent"}}'
```

AWS
re:Invent

aws

# Creating an automatic scaling policy
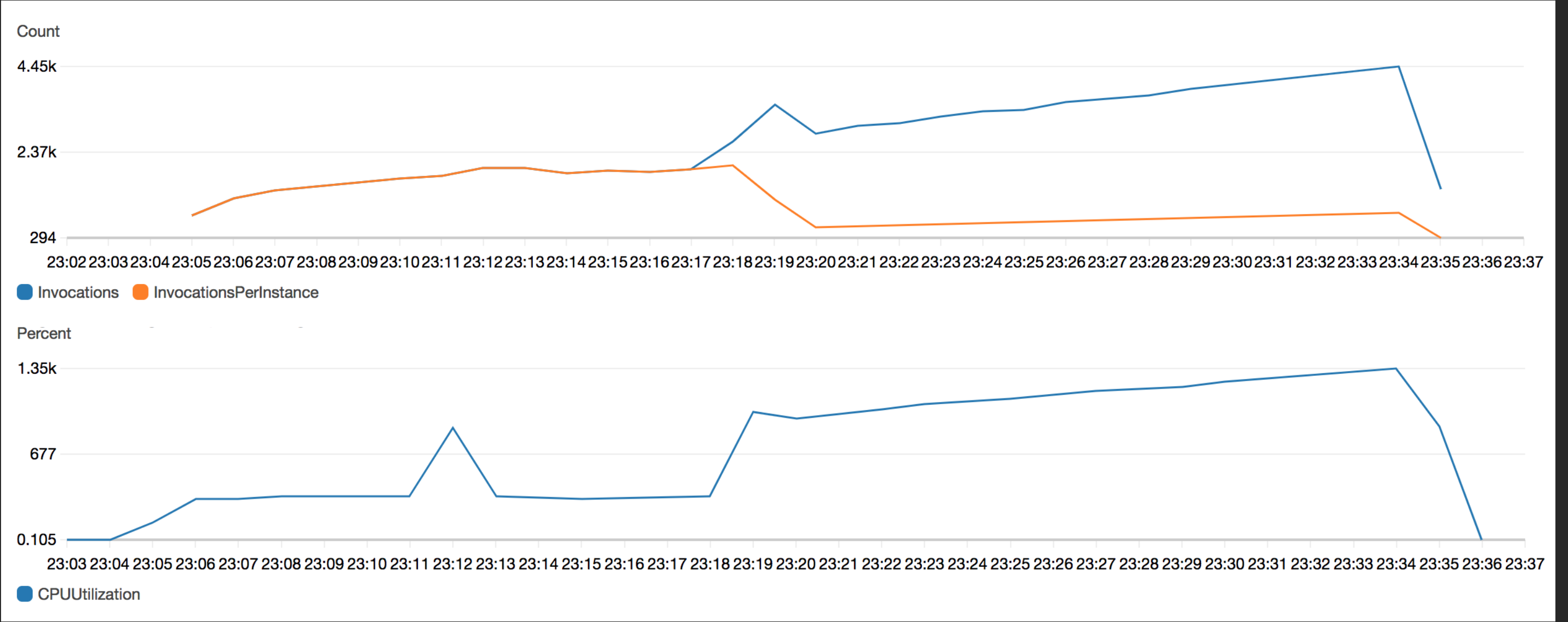
Variant

```
aws application-autoscaling register-scalable-target
    --service-namespace sagemaker
    --resource-id endpoint/my-endpoint/variant/model2
    --scalable-dimension sagemaker:variant:DesiredInstanceCount
    --min-capacity 2
    --max-capacity 5
```

Policy

```
aws application-autoscaling put-scaling-policy
    --policy-name model2-scaling
    --service-namespace sagemaker
    --resource-id endpoint/my-endpoint/variant/model2
    --scalable-dimension sagemaker:variant:DesiredInstanceCount
    --policy-type TargetTrackingScaling
    --target-tracking-scaling-policy-configuration
    '{"TargetValue": 50,
        "CustomizedMetricSpecification":
        {"MetricName": "CPUUtilization",
         "Namespace": "/aws/sagemaker/Endpoints",
         "Dimensions":
            [{"Name": "EndpointName", "Value": "my-endpoint"},
             {"Name": "VariantName","Value": "model2"}],
        "Statistic": "Average",
        "Unit": "Percent"}}'
```

AWS re:Invent

aws

# Scale by utilization

Intuit

aws

# Agenda

Who is Intuit?

Data lake functional architecture

Model development workflow

Key benefits of Amazon SageMaker

Standardizing model development for speed

Demo

aws

# Overall data lake architecture

# Model development workflow

# Key benefits of Amazon SageMaker

**Functional**

- Many algorithms supported with more being added constantly
- Custom algorithm supporting using docker
- Highly customizable
- Out of the box model parameter tuning

**Security**

- Integration with AWS Identity and Access Management (IAM) and AWS Key Management Service (AWS KMS)
- Good model for authentication and authorization

**Scalability of Amazon cloud**

# Standardizing model development for speed

## Standardized notebooks

- Standardized the security model around Amazon SageMaker
- Added functional integrations with Hive and our data marts

## Training

- Python library to make docker images that work with Amazon SageMaker training
- We use out-of-the-box Amazon SageMaker training

## Hosting

- Python library to make Docker images that work with Amazon SageMaker hosting
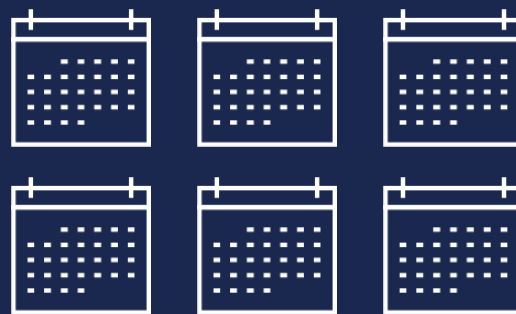- Integration with our internal Amazon API Gateway/Services Gateway

## Model deployment tool

- Tracking trained model versions and taking it through MDLC

intuit®

Deployment time down by 90%
with Amazon SageMaker

**6 MONTHS**

**< 1 WEEK**

# Demo



```
EXPLORER

▶ OPEN EDITORS
◢ UNTITLED (WORKSPACE)
  ◢ model-template                    ●
    ▶ dist
    ◢ model
        __init__.py
        main.py
    ▶ model.egg-info
      .gitignore
      api_creds                       U
    ⓘ README.md
      setup.py
  ▶ de-ops-model                      ●
```

Tabs: combined_dataset.csv | learning.py | README.md | Preview README.md | setup.py ✕

```python
 2
 3  setup(
 4      # Package information
 5      name='model',
 6      version='0.0.1',
 7
 8      # Package data
 9      packages=find_packages(),
10      include_package_data=True,
11
12      # Insert dependencies list here
13      install_requires=[
14          'pandas',
15          'tensorflow',
16          'tables',
17          'putz'
18      ],
19      entry_points={
20          'setuptools_docker.predict': [
21              'my_prediction_entrypoint = model.main:my_prediction_function',
22          ],
23          'setuptools_docker.train': [
24              'my_prediction_entrypoint = model.main:my_train_function',
25          ]
26      }
27  )
28
```
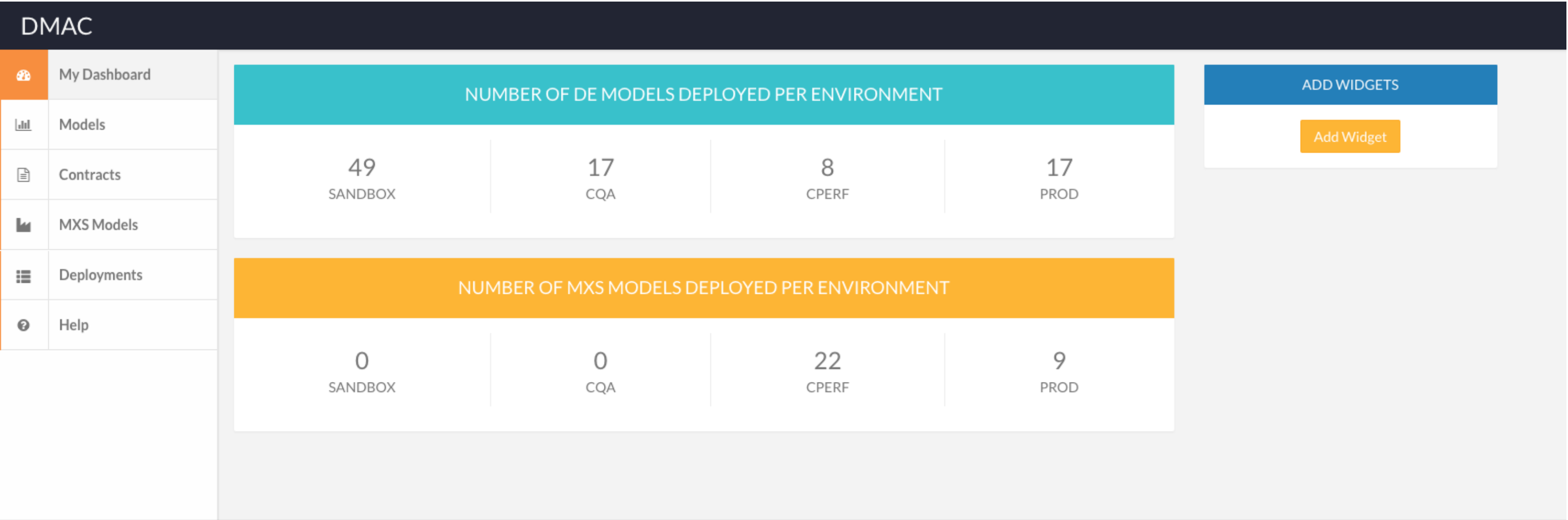
AWS re:Invent

aws

# Decision Model Automation Console (DMAC)

**DMAC**

- My Dashboard
- Models
- Contracts
- MXS Models
- Deployments
- Help

**NUMBER OF DE MODELS DEPLOYED PER ENVIRONMENT**

| 49 | 17 | 8 | 17 |
|---|---|---|---|
| SANDBOX | CQA | CPERF | PROD |

**NUMBER OF MXS MODELS DEPLOYED PER ENVIRONMENT**

| 0 | 0 | 22 | 9 |
|---|---|---|---|
| SANDBOX | CQA | CPERF | PROD |

**ADD WIDGETS**

Add Widget

AWS re:Invent

aws

# DMAC: post training



Enter Inference Service Parameters

GitHub Repo URL

https://github.intuit.com/pprabhu/model-test

GitHub Repo Tag

master

S3 Model Data Path

s3://test-model-data
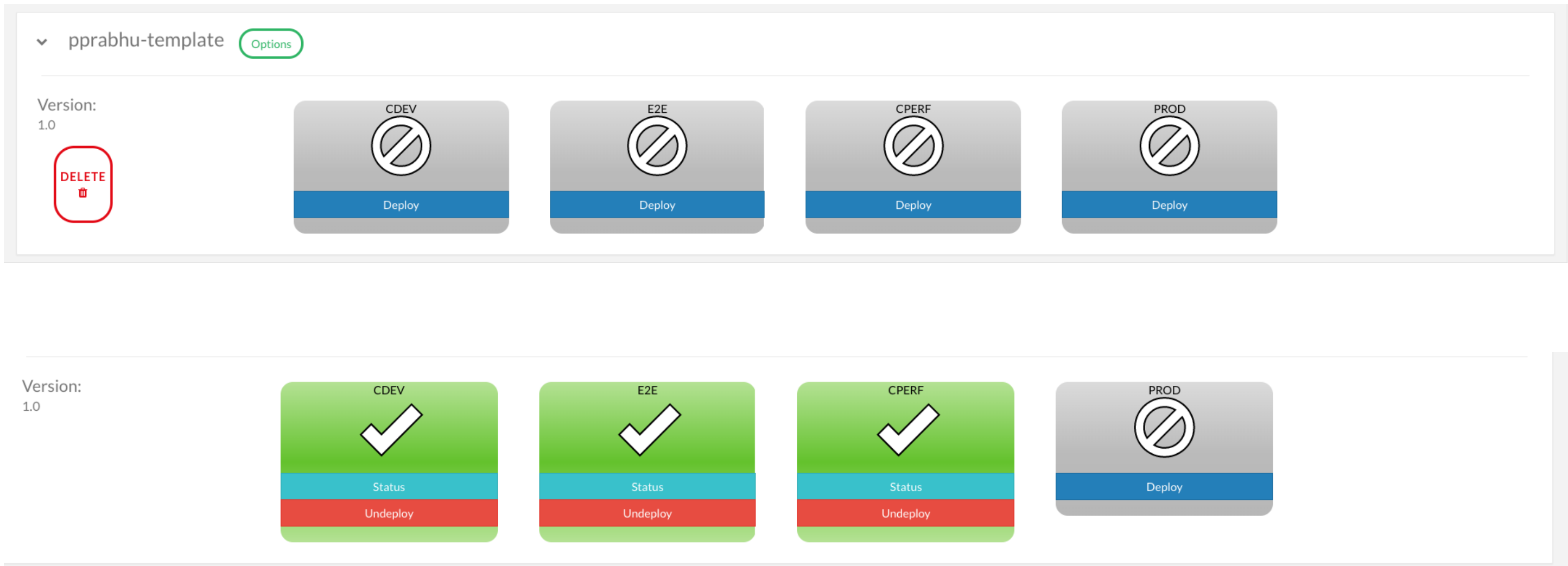
Choose Model Type

Custom Python 3.6 ▼

BACK    CANCEL    SUBMIT

# Model deployment workflow

# Conclusion

- Amazon SageMaker is a versatile platform to build, train, and deploy machine learning models at scale.

- Customers like Intuit are benefiting from integrating Amazon SageMaker into their data science workflows.

- Explore Amazon SageMaker (free tier eligible*) and build models of your own.

*https://aws.amazon.com/free/

# Thank you!

David Arpin
Prasad Prabhu

aws

Please complete the session survey in the mobile app.