AWS
re:Invent

AIM326-R

# Implementing ML workflows with Kubernetes and Amazon SageMaker

**David Ping**
Principal ML Solutions Architect
Amazon Web Services

**Aditya Bindal**
Senior Product Manager
Amazon Web Services

**Suhas Guruprasad**
ML Lead
Zalando

aws
re:Invent

aws

# Agenda

- Machine learning is a hard problem
- Fully managed machine learning with Amazon SageMaker
- Kubernetes architecture
- Build ML workflows with Amazon SageMaker and Kubernetes
- Machine learning at Zalando
- Demo

Your organization has adopted  to modernize technology

Need to transform business with machine learning
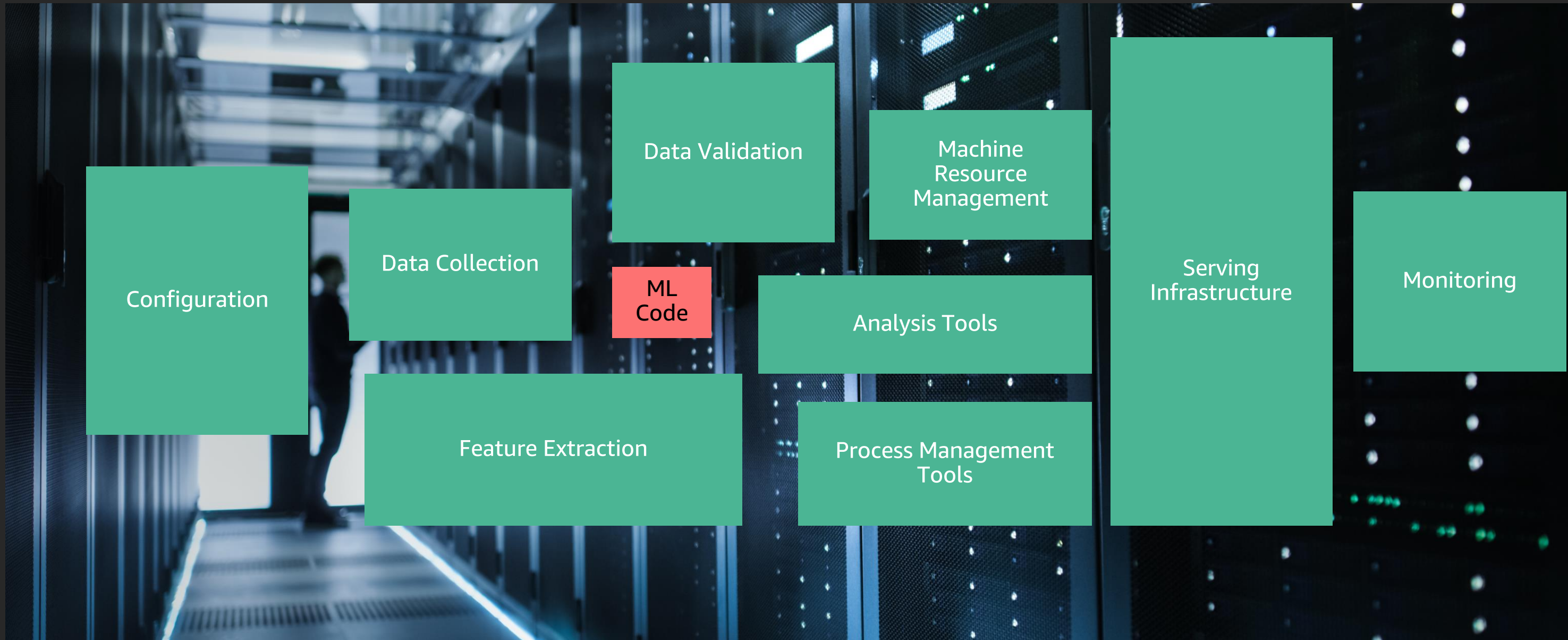
Data scientist

K8s DevOps engineer

- Focuses on data science, business outcome and speed to market

- Wants minimum dependency on the DevOps team for experimentation and model development
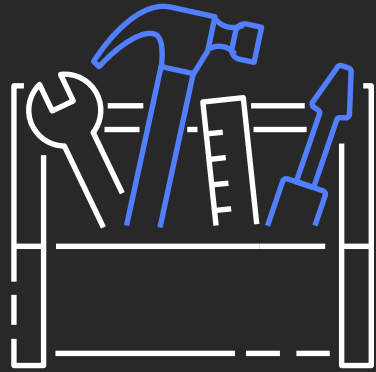
- No or limited K8s and infrastructure knowledge

- Wants to leverage existing K8s investment and best practices

- Wants to manage using familiar K8s construct and syntax

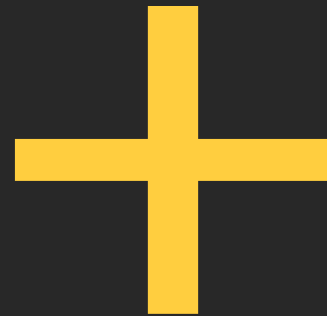- Limited ML knowledge & engineering experience with ML workloads

# Machine learning is hard



Configuration

Data Collection

Data Validation

ML Code

Machine Resource Management

Feature Extraction

Analysis Tools

Process Management Tools

Serving Infrastructure
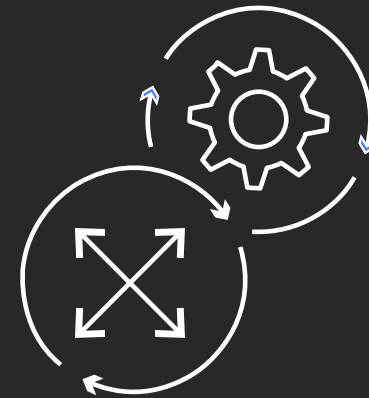
Monitoring

Source: Sculley et al: Hidden Technical Debt in Machine Learning Systems

Do-it-yourself          Opportunity
                        for both?          Managed service

# Fully managed machine learning with Amazon SageMaker

aws

# Amazon SageMaker is a fully managed service that covers the entire machine learning workflow
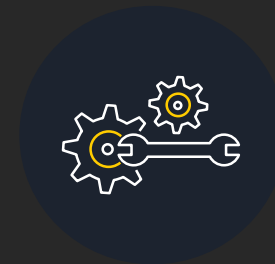
Jupyter notebook instances

High-performance algorithms

Large-scale training
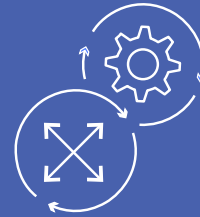
Optimization

One-click deployment

Fully managed with auto-scaling

# Built on modern application architecture

## Containers

- Built-in algorithms for training and hosting
- Managed ML containers for training and hosting
- Bring-your-own containers for training and hosting

## Serverless

- No server to manage
- API driven
- Pay by usage
- Built-in monitoring and logging

## Cloud Native

- Multi-AZ deployment
- Scale with powerful compute resources
- IAM, VPC, Encryption
- Resource on-demand or SPOT, Amazon Elastic Inference, Multi-Model Endpoint

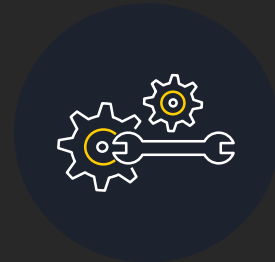# Minimized dependency on DevOps engineering

Data scientist

aws

Jupyter notebook instances

High-performance algorithms

Large-scale training

Optimization

One-click deployment

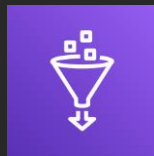Fully managed with auto-scaling

# Run Amazon SageMaker in a pipeline
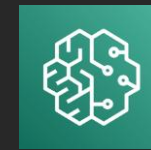
DevOps engineer
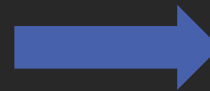
Workflow orchestration      AWS Step Functions
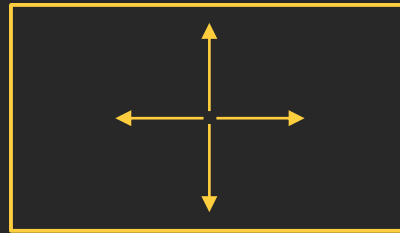
AWS Glue

**Data Pipeline**

Amazon SageMaker

**ML Pipeline**

# Kubernetes architecture

aws

# What is Kubernetes (aka K8s)?
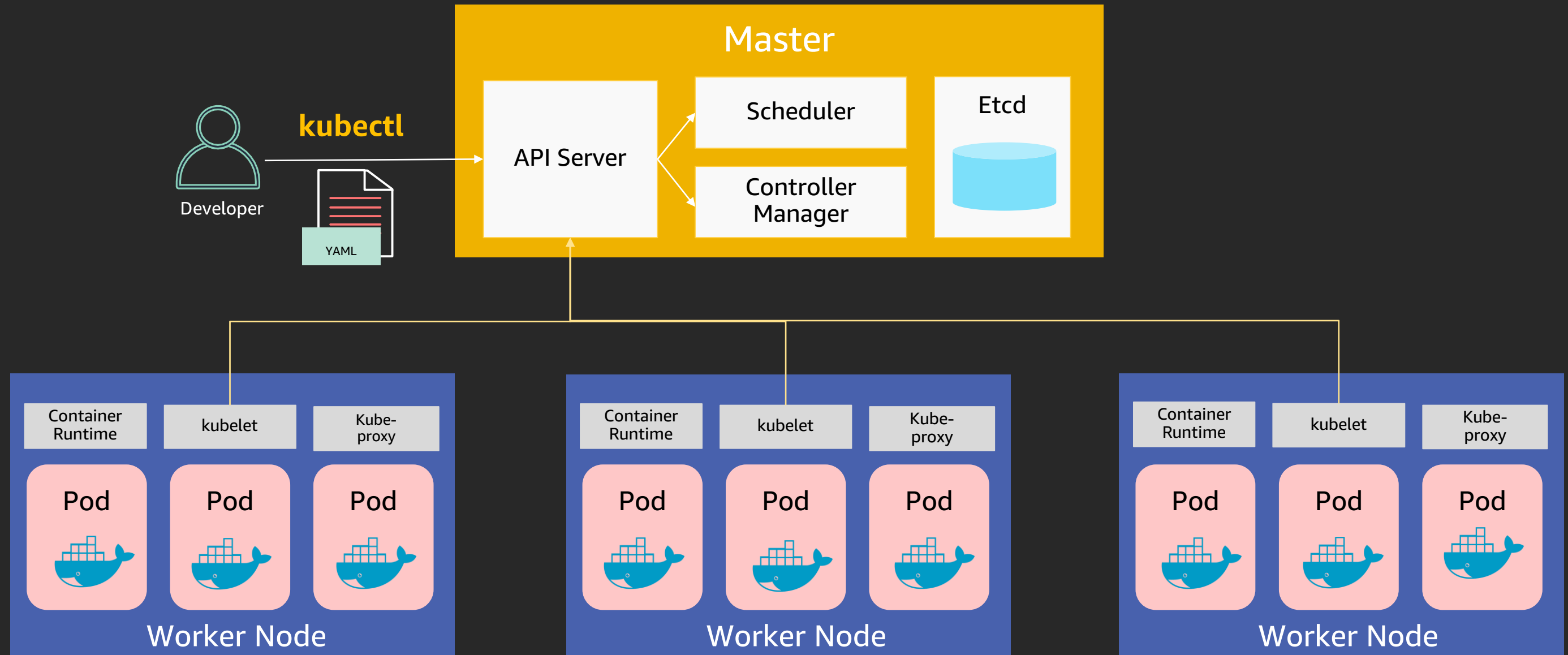
**Open-source container management platform**

**Helps you run containers at scale**

**Gives you primitives for building modern applications**

# Kubernetes architecture

# Building ML workflows with Kubernetes and Amazon SageMaker

**Using Kubernetes for ML is hard to manage and scale**

Build and manage services within Kubernetes cluster for ML

+

Make disparate open-source libraries and frameworks work together in a secure and scalable way

+

Requires time and expertise from infrastructure, data science, and development teams

=

Need an easier way to use Kubernetes for ML

# Introducing Amazon SageMaker Operators for Kubernetes

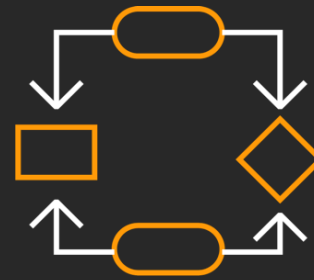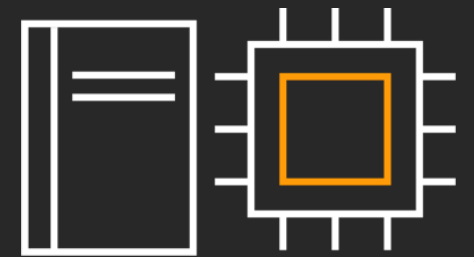## Kubernetes customers can now train, tune, and deploy models in Amazon SageMaker

Train, tune, and deploy models in Amazon SageMaker
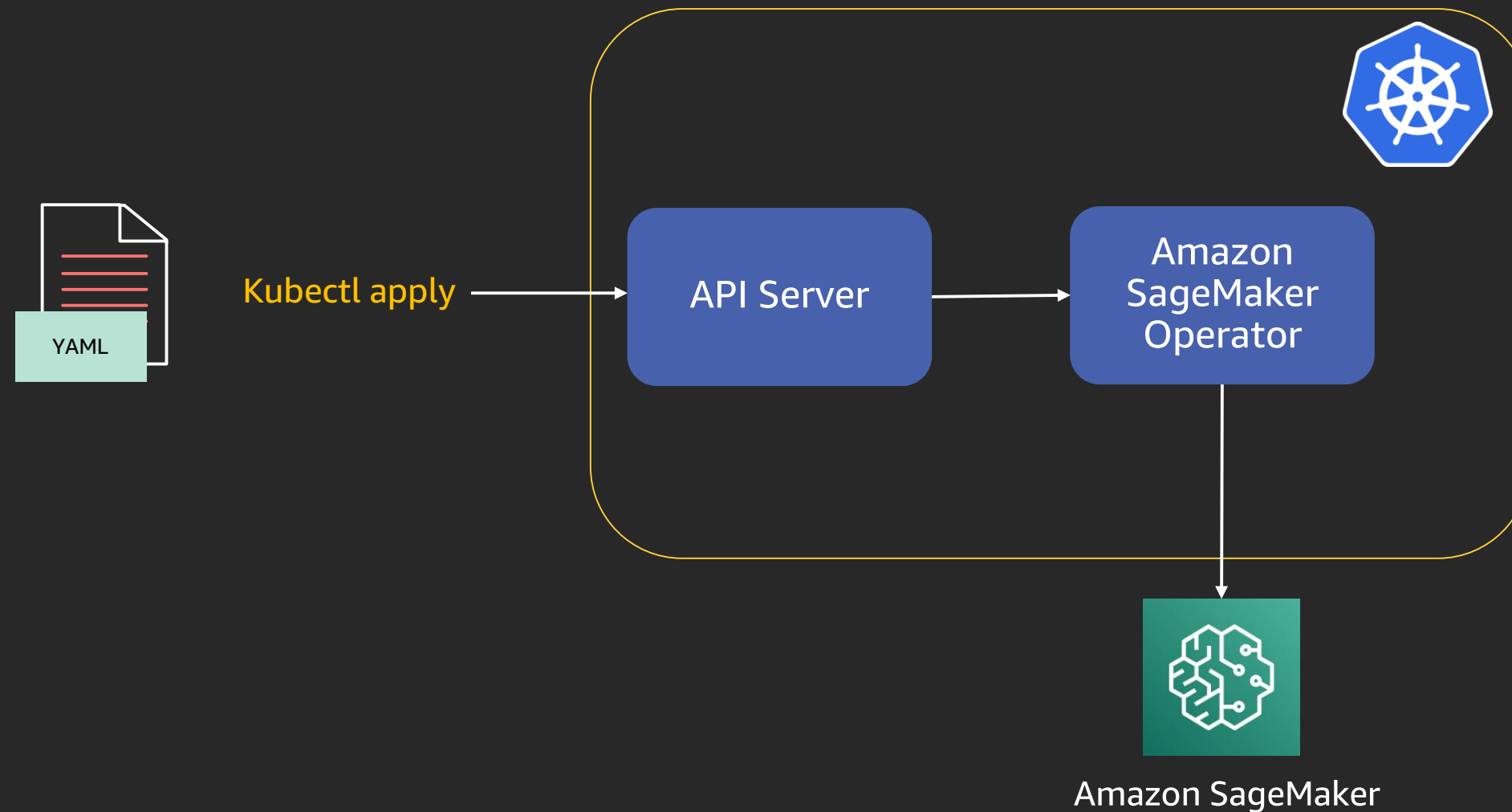
Orchestrate ML workloads from your Kubernetes environments

Create pipelines and workflows in Kubernetes

Fully managed infrastructure in Amazon SageMaker

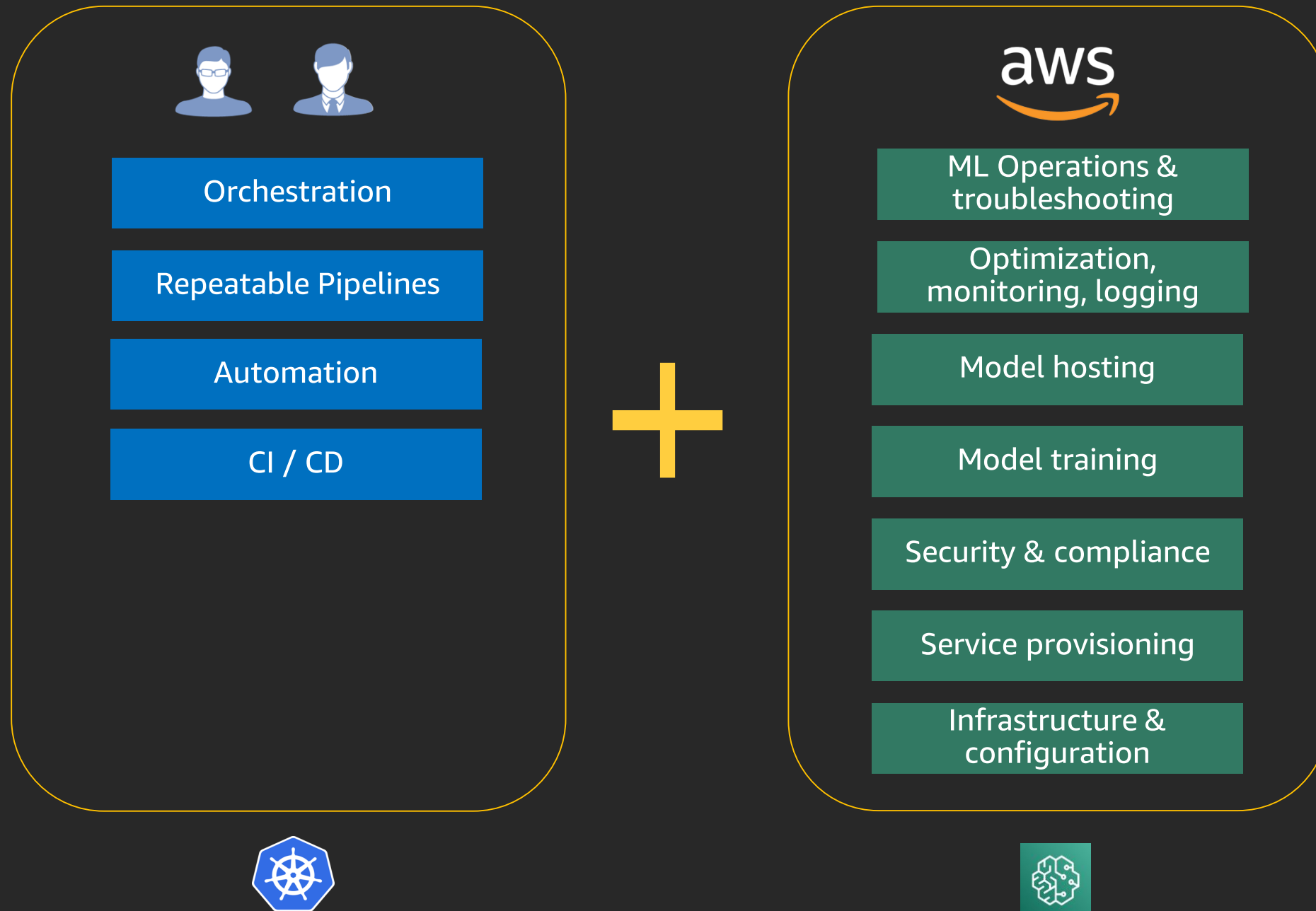# Under the hood – Amazon SageMaker and Kubernetes

YAML

**Kubectl apply** →

API Server → Amazon SageMaker Operator

Amazon SageMaker

**Key Features**

- Amazon SageMaker Operators for training, tuning, inference

- Natively interact with Amazon SageMaker jobs using K8s tools (e.g., get pods, describe)

- Stream and view logs from Amazon SageMaker in K8s

- Helm Charts to assist with setup and spec creation

# Why together?

Orchestration

Repeatable Pipelines

Automation

CI / CD

**+**

aws

ML Operations & troubleshooting

Optimization, monitoring, logging

Model hosting

Model training

Security & compliance

Service provisioning

Infrastructure & configuration

# Machine learning at Zalando

AWS
re:Invent

aws

Zalando.
The starting point for fashion.

zalando

# ZALANDO AT A GLANCE

**~ 5.4** billion EUR

**revenue 2018**

**~ 14,000** employees in Europe

**> 80%** of visits via mobile devices

**> 300 million** visits per month

**> 29 million** active customers

**> 450,000** product choices

**> 2,000** brands

**17** countries

zalando

# "A sustainable fashion platform with a net positive impact for people and planet"

zalando

# We leverage Machine Learning across the platform to deliver a better customer experience.

zalando

# MACHINE LEARNING AT ZALANDO

## Get the Look

Perfectly combined by fashion experts. Get fresh outfit ideas for the new season and every occasion.

**All Styles**    Classic    Casual    Trendy    Party

## Inspired by you

zalando

# MACHINE LEARNING AT ZALANDO

# MACHINE LEARNING AT ZALANDO

# MACHINE LEARNING AT ZALANDO

zalando

**But what constitutes an ML project and how is it different from traditional software?**

zalando

# MICROSERVICES AT ZALANDO

**CODE**

**and commits to GitHub**

zalando

# MICROSERVICES AT ZALANDO

**CODE**

and commits to GitHub

**CONTINUOUS DEPLOYMENT**

builds triggered by pull requests

zalando

# MICROSERVICES AT ZALANDO

**CODE**
and commits to GitHub

**CONTINUOUS DEPLOYMENT**
builds triggered by pull requests

**KUBERNETES**
Software runs on one of our >140 Kubernetes clusters on AWS

zalando

# CONTINUOUS DEPLOYMENT

# THE MACHINE LEARNING JOURNEY



**ANALYSIS**
**Exploratory Data**
**Analysis**

zalando

# THE MACHINE LEARNING JOURNEY

**ANALYSIS**
Exploratory Data
Analysis

**TRAINING**
Features, Training &
Evaluating models

zalando

# THE MACHINE LEARNING JOURNEY

**ANALYSIS**

**Exploratory Data Analysis**

**TRAINING**

**Features, Training & Evaluating models**

**OBSERVATIONS**

**Run A/B tests, see KPIs & adjust**

zalando

# Building a central ML platform: What are the challenges?

zalando

# CHALLENGES

**SPEED**
Research to
production time

zalando

# CHALLENGES

**SPEED**

Research to
production time

**SAFETY**

Understanding, testing,
reproducibility, monitoring

zalando

# CHALLENGES

**SPEED**

**Research to production time**

**SAFETY**

**Understanding, testing, reproducibility, monitoring**

**SCALE**

**Run thousands of experiments**

zalando

# CHALLENGES



**COST-EFFICIENCY**

**Experiments must be cheap at scale**

zalando

# CHALLENGES

**COST-EFFICIENCY**
Experiments must be
cheap at scale

**COLLABORATION**
Across multiple job
families

# CHALLENGES

**COST-EFFICIENCY**
Experiments must be cheap at scale

**COLLABORATION**
Across multiple job families

**SATISFACTION**
Job satisfaction and happiness

# Relooking at the ML journey to define the ML pipeline for the ML platform
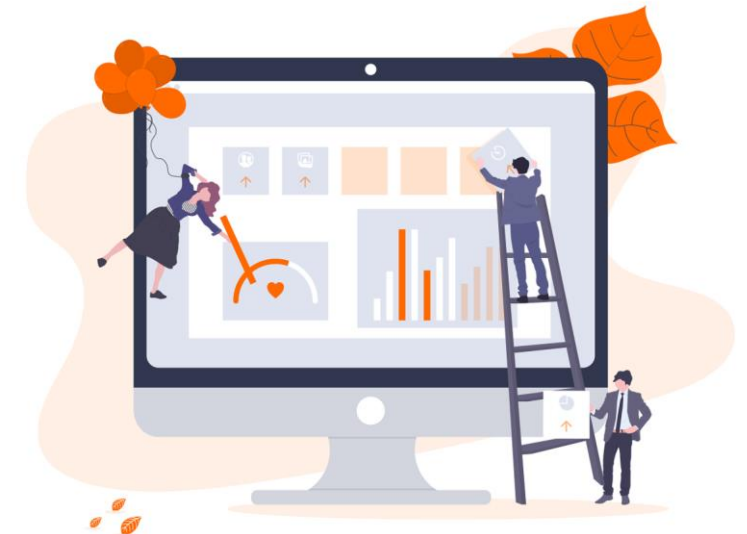
zalando

# THE MACHINE LEARNING JOURNEY

**ANALYSIS**
**Exploratory Data Analysis**

**TRAINING**
**Features, Training & Evaluating models**

**OBSERVATIONS**
**Run A/B tests, see KPIs & adjust**

zalando

# MACHINE LEARNING PIPELINES

**FUNDAMENTAL**

**ML should be about pipelines,**

**not just data or models**

zalando

# Using Amazon SageMaker for ML pipelines at Zalando

zalando

# WHY AMAZON SAGEMAKER AT ZALANDO

**TRAINING JOBS**

Training jobs repr.
core of a pipeline.

zalando

# WHY AMAZON SAGEMAKER AT ZALANDO

## TRAINING JOBS

Training jobs repr. core of a pipeline.

## ALGORITHMS

Built-in and bring-your-own docker images.

zalando

# WHY AMAZON SAGEMAKER AT ZALANDO

**TRAINING JOBS**

**Training jobs repr. core of a pipeline.**

**ALGORITHMS**

**Built-in and bring-your-own docker images.**

**INTEGRATED**

**With AWS offerings.**

zalando

# WHY AMAZON SAGEMAKER AT ZALANDO

**TRAINING JOBS**

Training jobs repr. core of a pipeline.

**ALGORITHMS**

Built-in and bring-your-own docker images.

**INTEGRATED**

With AWS offerings.

**"SERVERLESS"**

Training instances are managed, jobs easy to distribute.

zalando

# WHY AMAZON SAGEMAKER AT ZALANDO

## TRAINING JOBS
Training jobs repr. core of a pipeline.

## ALGORITHMS
Built-in and bring-your-own docker images.

## INTEGRATED
With AWS offerings.

## "SERVERLESS"
Training instances are managed, jobs easy to distribute.

## BATCH & APIs
In the same package. Easy to reason and deploy.

zalando

# WHY AMAZON SAGEMAKER AT ZALANDO

**TRAINING JOBS**

Training jobs repr. core of a pipeline.

**ALGORITHMS**

Built-in and bring-your-own Docker images.

**INTEGRATED**

With AWS offerings.

**"SERVERLESS"**

Training instances are managed, jobs are easy to distribute.

**BATCH & APIs**

In the same package. Easy to reason and deploy.

**SCRIPT MODE**

Easy to write quick TensorFlow scripts and tie in to pipelines

zalando

# ZALANDO MACHINE LEARNING PLATFORM



**ZFLOW**

Zalando's pythonic interface to
ML pipelines

**STEP FUNCTIONS**

To orchestrate pipelines

**SAGEMAKER**

The Core Engine

56

zalando

# ZALANDO MACHINE LEARNING PLATFORM

```
training_stage = training_job(
    name="...",
    input_data_configs="...",
    hyperparameters="..."
)

batch_transform_stage =
batch_transform_job(
    name="...",
    s3_input_path="...",
    s3_output_path="..."
)
```

zalando

# ZALANDO MACHINE LEARNING PLATFORM

```
pipeline = PipelineBuilder("my_pipeline")

pipeline
    .add_stage(data_processing_stage)
    .add_stage(training_stage)
    .add_stage(batch_transform_stage)
```

zalando

# ZALANDO MACHINE LEARNING PLATFORM

≡ **DEPLOYMENT UNITS**   📁 **ZFLOW-EXAMPLE**   ⌥ **PR-26-22**

| ZFLOW-EXAMPLE | ⚙ UPLOAD-ARTIFACTS-TO-S3 | ☁ DEPLOY-PIPELINE |
|---|---|---|

**GENERAL**   **DEPLOYMENT STATUS**                                          ⏱ Ran for 3 m, 12 s

✅ CloudFormation **ChurnPredictionPipelineStack** console ⧉                                ⊖

   ✅ AWS::Lambda::Function **CheckDatabricksJobStatus111D2B43**

   ✅ AWS::IAM::Role **ChurnPredictionPipelineEventsRoleC566F7A2**

   ✅ AWS::IAM::Policy **ChurnPredictionPipelineEventsRoleDefaultPolicyC77F2FF9**

   ✅ AWS::StepFunctions::StateMachine **ChurnPredictionPipelineF3E3A849**

   ✅ AWS::IAM::Role **ChurnPredictionPipelineRole476B0573**

   ✅ AWS::IAM::Policy **ChurnPredictionPipelineRoleDefaultPolicyC8A22D3C**

   ✅ AWS::KMS::Key **ChurnPredictionPipelineStackKMSKeyA32054EA**

59

zalando

# UPCOMING



**ZFLOW**

**K8S operators**

Interfacing via kubectl and other software deployments

**STEP FUNCTIONS**

To orchestrate pipelines

**SAGEMAKER**

The Core Engine

zalando

# ZALANDO MACHINE LEARNING PLATFORM

```yaml
apiVersion: sagemaker.aws.amazon.com/v1
kind: TrainingJob
metadata
    name: kmeans-mnist
spec
    algorithmSpecification
        trainingImage: .../kmeans:1
        trainingInputMode: File
        hyperParameters: ...
        inputDataConfig:
    resourceConfig
        instanceType: ml.c4.8xlarge
    ...
```

zalando

# OVERCOMING THE CHALLENGES

**SPEED**

**Notebooks + consistent interfaces**

zalando

# OVERCOMING THE CHALLENGES
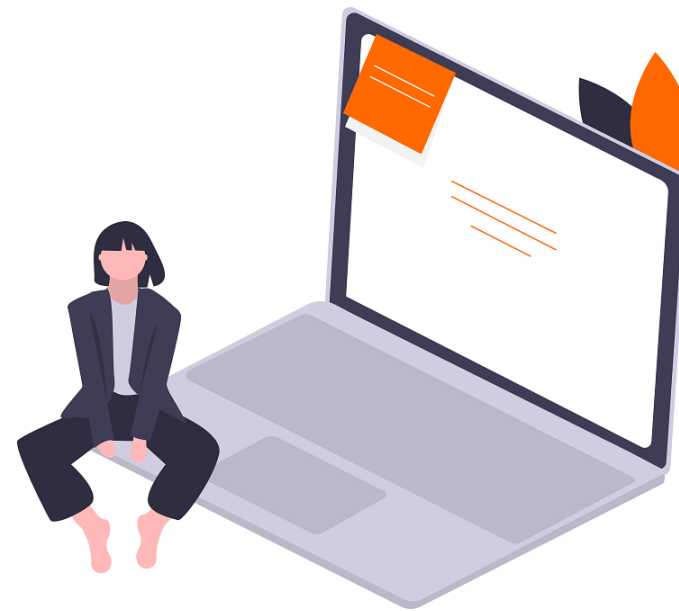
**SPEED**

**Notebooks + consistent**

**interfaces**

**SAFETY**

**Amazon**

**SageMaker**

**training/scoring**

zalando

# OVERCOMING THE CHALLENGES

**SPEED**
**Notebooks + consistent interfaces**

**SAFETY**
**Amazon SageMaker training/scoring**

**SCALE**
**Amazon SageMaker training/scoring**

zalando

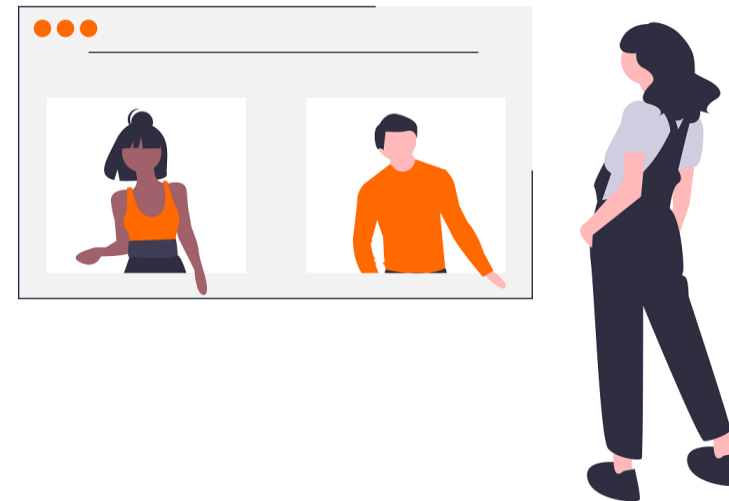# OVERCOMING THE CHALLENGES

**COST-EFFICIENCY**

**Amazon SageMaker**

**training/scoring**

zalando

# OVERCOMING THE CHALLENGES

**COST-EFFICIENCY**
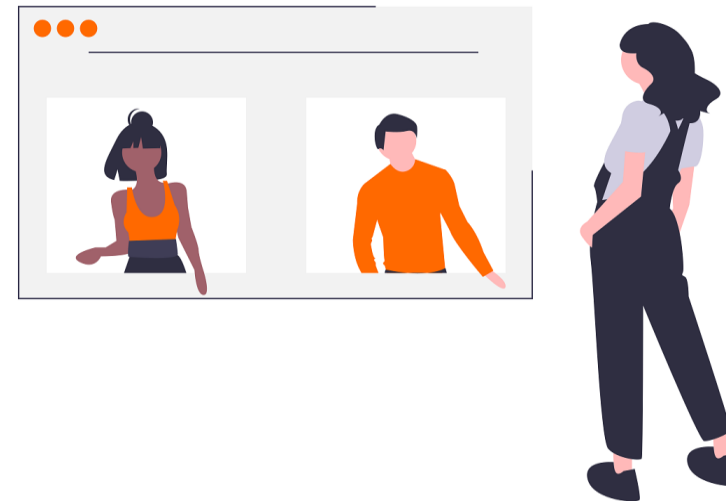**Amazon SageMaker**
**training/scoring**

**COLLABORATION**
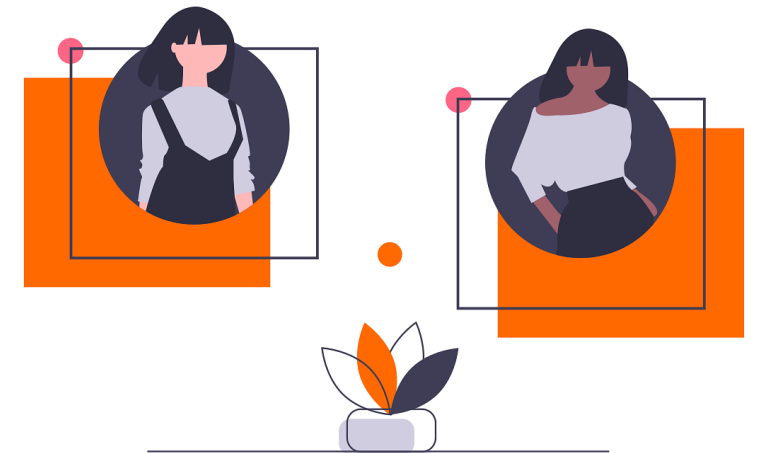**Consistent interfaces**

zalando

# OVERCOMING THE CHALLENGES

**COST-EFFICIENCY**

Amazon SageMaker training/scoring

**COLLABORATION**

Consistent interfaces

**SATISFACTION**

Central team support

zalando

**Machine learning at scale offers exciting new opportunities for us to serve our customers better.**

zalando

# Demo

aws

# In summary

- Machine learning is a hard problem

- Amazon SageMaker simplifies ML with modern application architecture

- Kubernetes is great for large-scale container orchestration

- Amazon SageMaker and Kubernetes can provide greater benefits when combined

# Thank you!

aws re:Invent

aws

Please complete the session survey in the mobile app.