



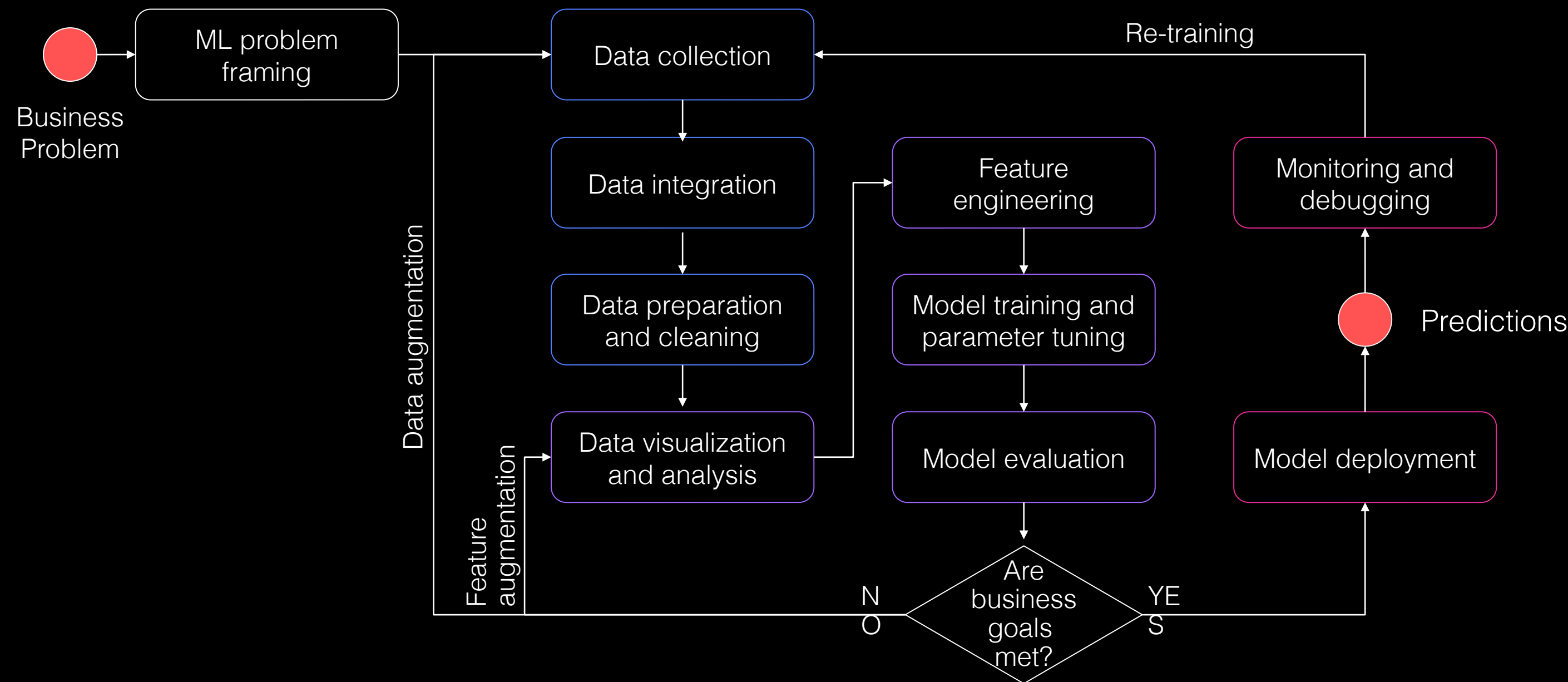
SUMMIT  
Switzerland

# Build, Train and Deploy Machine Learning Models on Amazon SageMaker

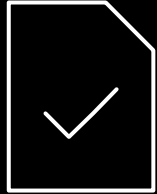
Julien Simon  
Global Evangelist, AI & Machine Learning  
Amazon Web Services  
@julsimon

Stéphane Cheikh  
Director, Portfolio Evolution using Artificial Intelligence  
SITA

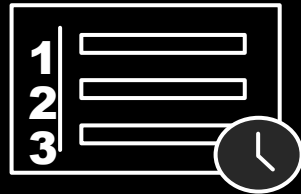
# Machine learning cycle



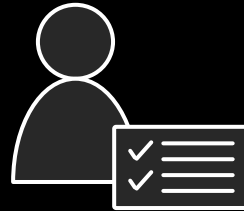
# Amazon SageMaker



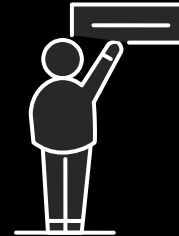
Collect and prepare  
training data



Choose and  
optimize your  
ML algorithm



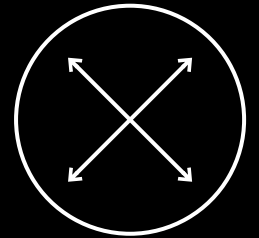
Set up and  
manage  
environments  
for training



Train and  
tune ML models



Deploy models  
in production



Scale and manage  
the production  
environment

Same service and APIs from experimentation to production

intuit



tinder



CONVOY

SIEMENS



DOW JONES



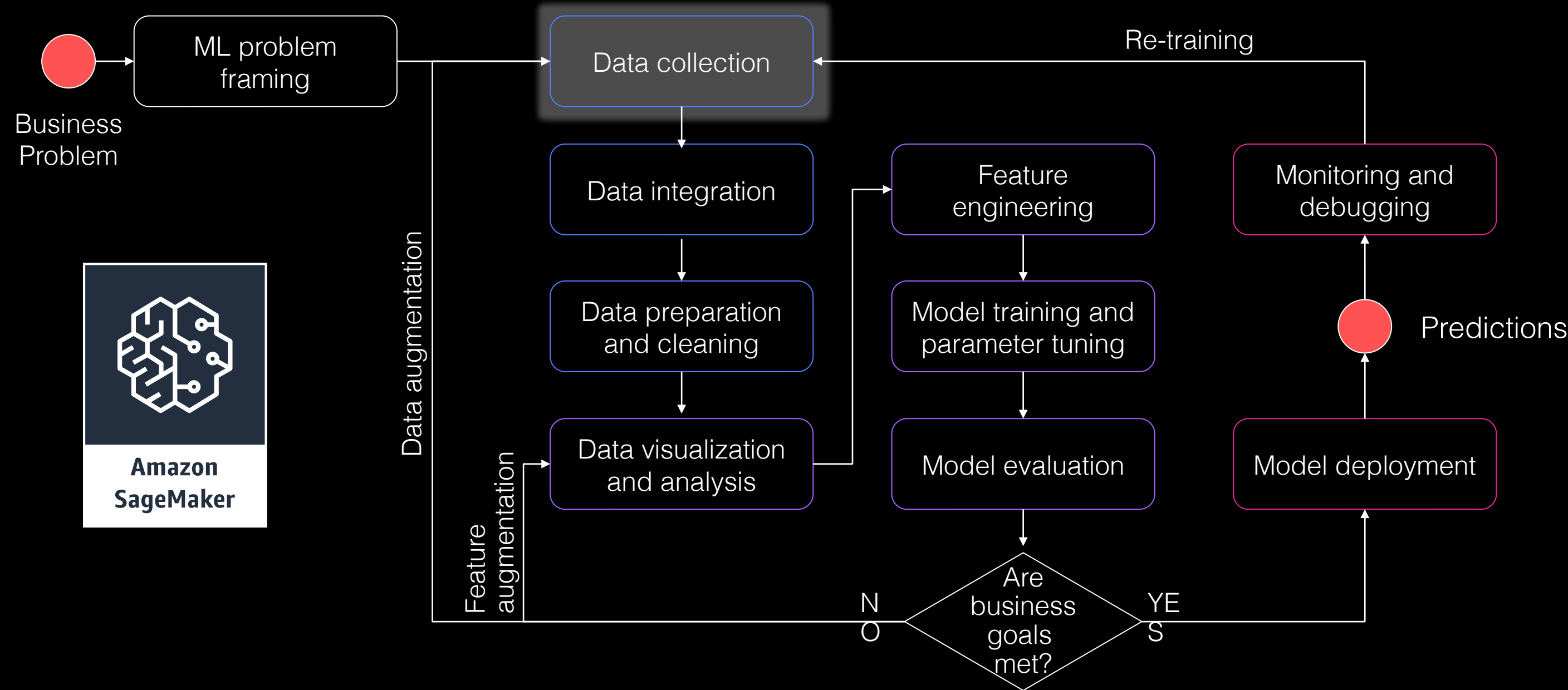
SONY



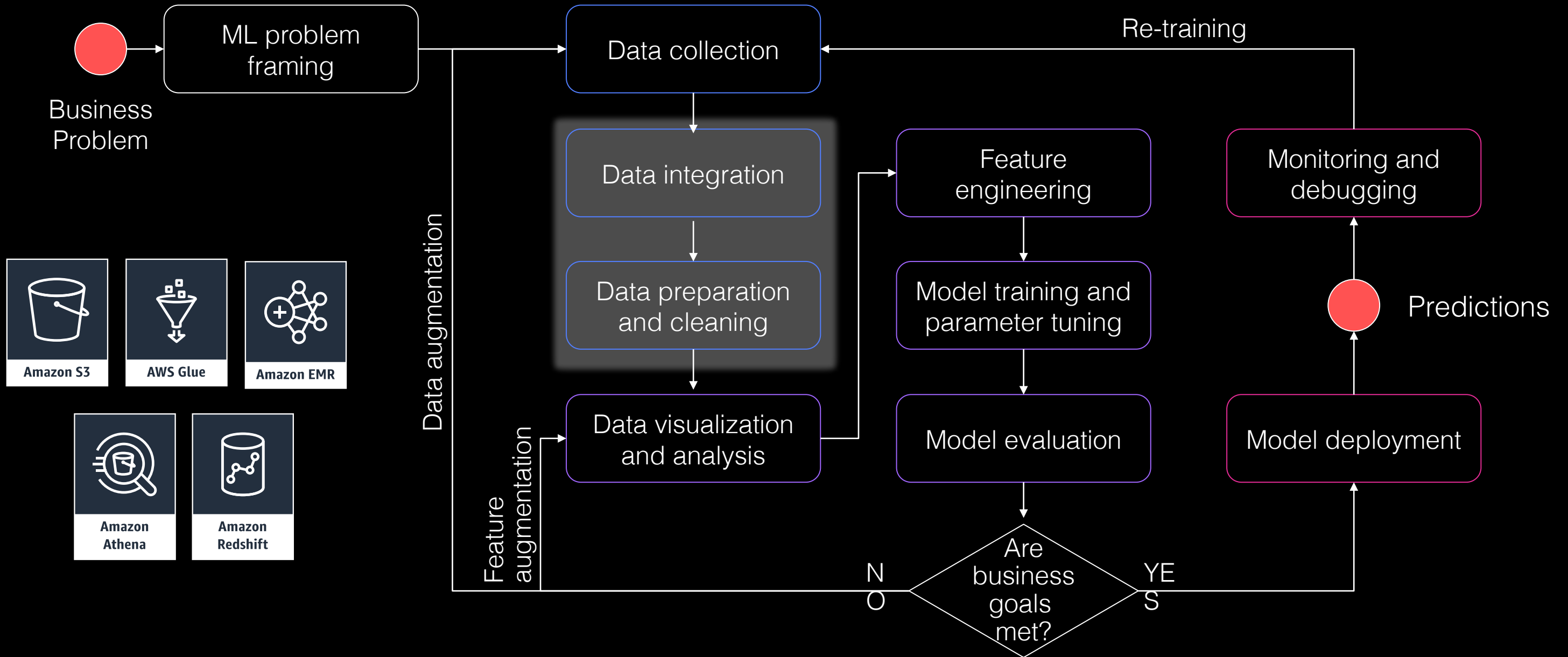
aws SUMMIT

© 2019, Amazon Web Services, Inc. or its affiliates. All rights reserved.

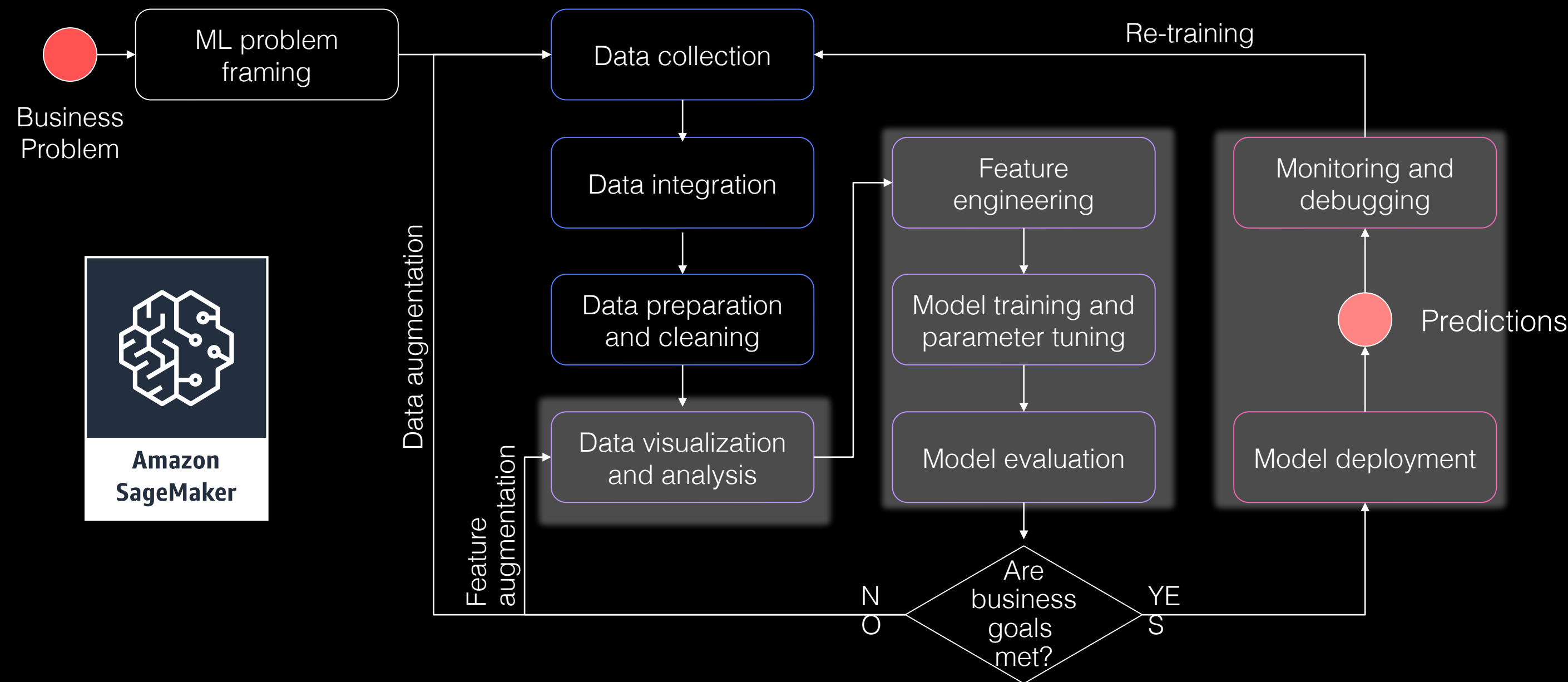
# Build your dataset



# Prepare your dataset for Machine Learning



# Build, train and deploy models using SageMaker

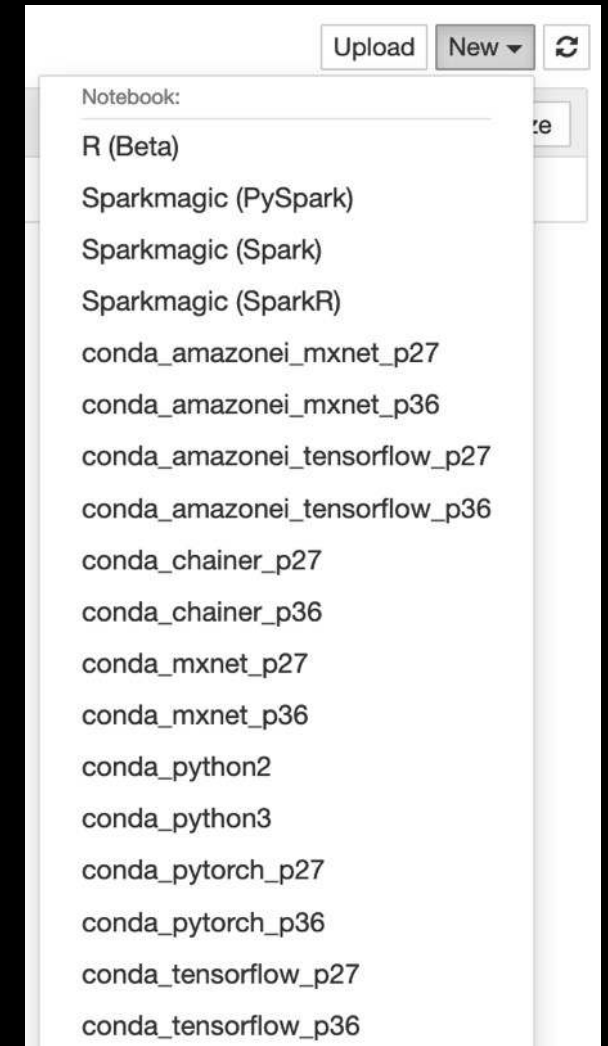


# Building models with Amazon SageMaker



# Notebook instances

- *Fully managed EC2 instances, from T2 to P3*
  - G4 and R5 now available for inference – **NEW!**
- Pre-installed with **Jupyter** and **Conda** environments
  - Python 2.7 & 3.6
  - Open-source libraries (TensorFlow, Apache MXNet, etc.)
  - Beta support for R – **NEW!**
  - Amazon Elastic Inference for cost-effective GPU acceleration
- Lifecycle configurations
- VPC, encryption, etc.
- Get to work in minutes, **zero setup**



# The Amazon SageMaker API

- Python SDK **orchestrating** all Amazon SageMaker activity
  - High-level objects for **algorithm selection**, **training**, **deploying**, **automatic model tuning**, etc.  
<https://github.com/aws/sagemaker-python-sdk>
  - **Spark SDK** (Python & Scala)  
<https://github.com/aws/sagemaker-spark/tree/master/sagemaker-spark-sdk>
- AWS SDK
  - Service-level APIs for **scripting** and **automation**
  - CLI: *'aws sagemaker'*
  - Language SDKs: boto3, etc.

NEW  
!

Ground Truth

Training data

Amazon S3

Model artifacts

Amazon S3  
Amazon EFS  
Amazon FSx for Lustre

Client application

Inference response

Inference request

Inference endpoint

Amazon SageMaker

Amazon ECR



Inference code



Helper code

Model Hosting



Training code



Helper code

Model Training (on demand or spot)

NEW  
!



Inference code



Training code

# Model options



Training code

**AWS Marketplace  
for Machine  
Learning:  
250+ off-the-shelf  
algorithms and models**

Factorization Machines  
Linear Learner  
Principal Component Analysis  
K-Means Clustering  
XGBoost  
And more

Built-in Algorithms (17)

No ML coding required  
No infrastructure work required  
Distributed training  
Pipe mode



Built-in Frameworks

Bring your own code: Script mode  
Open-source containers  
No infrastructure work required  
Distributed training  
Pipe mode



Bring Your Own Container

Full control, run anything!  
R, C++, etc.  
No infrastructure work required

# Built-in algorithms

Orange: supervised, yellow: unsupervised

<b>Linear Learner:</b> Regression, classification	<b>Image Classification:</b> Deep learning (ResNet)
<b>Factorization Machines:</b> Regression, classification, recommendation	<b>Object Detection (SSD):</b> Deep learning (VGG or ResNet)
<b>K-Nearest Neighbors:</b> Non-parametric regression and classification	<b>Neural Topic Model:</b> Topic modeling
<b>XGBoost:</b> Regression, classification, ranking <a href="https://github.com/dmlc/xgboost">https://github.com/dmlc/xgboost</a>	<b>Latent Dirichlet Allocation:</b> Topic modeling (mostly)
<b>K-Means:</b> Clustering	<b>BlazingText:</b> GPU-based Word2Vec, and text classification
<b>Principal Component Analysis:</b> Dimensionality reduction	<b>Sequence to Sequence:</b> Machine translation, speech to text and more
<b>Random Cut Forest:</b> Anomaly detection	<b>DeepAR:</b> Time-series forecasting (RNN)
<b>Object2Vec:</b> General-purpose embedding	<b>IP Insights:</b> Usage patterns for IP addresses
<b>Semantic Segmentation:</b> Deep learning	

# Built-in frameworks: just add your code



- Built-in containers for **training** and **prediction**
  - Open-source, e.g., <https://github.com/aws/sagemaker-tensorflow-containers>
  - Build them, run them on your own machine, customize them, etc.
- **Local mode**: train and predict on your **notebook instance**, or on your **local machine**
- **Script mode**: migrate **existing code** to SageMaker with minimal changes

# TensorFlow on AWS

C5 instances (Intel Skylake)



Training ResNet-50 with the ImageNet dataset using our optimized build of TensorFlow 1.11 on a **c5.18xlarge** instance type is designed to be **11x faster** than training on the stock binaries

P3 instances (NVIDIA V100)

TensorFlow scaling efficiency with 256 GPUs

**65**

Stock version



**90**  
%

AWS-optimized version

# Apache MXNet: Deep learning for enterprise developers



## Start with off-the-shelf models

- Gluon CV, Gluon NLP, Gluon TS
- ONNX compatibility

## Fast and scalable training

- Keras-MXNet up to 2x faster than Keras-TensorFlow
- Near-linear scalability up to 256 GPUs
- Dynamic training

## Easy deployment

- Java and Scala APIs
- Model Server



# Demo:

## Image classification with Keras/TensorFlow

- + Script Mode
- + Managed Spot Training
- + Elastic Inference

<https://aws.amazon.com/blogs/machine-learning/train-and-deploy-keras-models-with-tensorflow-and-apache-mxnet-on-amazon-sagemaker/>

<https://gitlab.com/juliensimon/dlnotebooks/tree/master/keras/05-keras-blog-post>

# Getting started

<http://aws.amazon.com/free>

<https://ml.aws>

<https://aws.amazon.com/sagemaker>

<https://github.com/aws/sagemaker-python-sdk>

<https://github.com/aws/sagemaker-spark>

<https://github.com/aws-labs/amazon-sagemaker-examples>

<https://gitlab.com/juliensimon/dlnotebooks>

# Thank you!

Julien Simon  
Global Evangelist, AI & Machine Learning  
Amazon Web Services  
@julsimon

Stéphane Cheikh  
Director, Portfolio Evolution using Artificial Intelligence  
SITA



Please complete the  
session survey.