



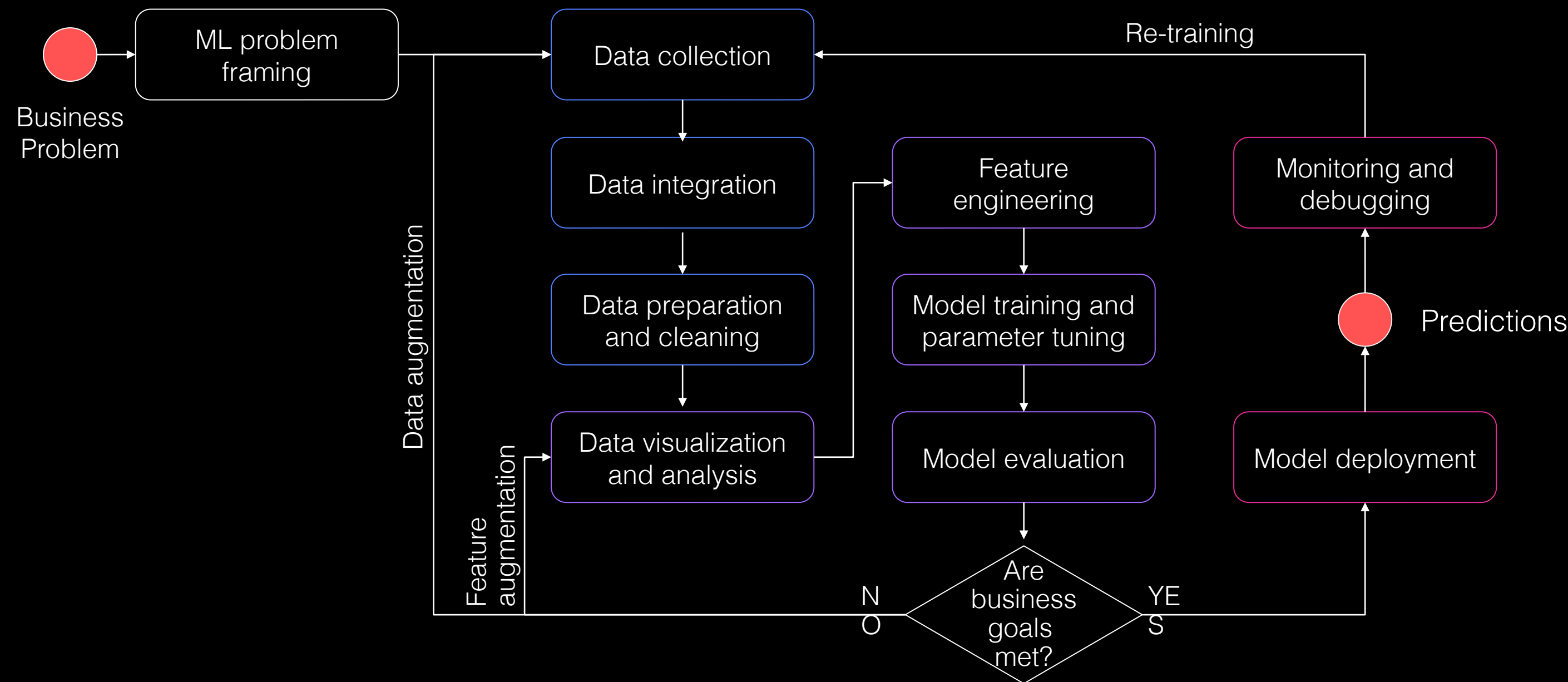
S U M M I T
Switzerland

Build, Train and Deploy Machine Learning Models on Amazon SageMaker

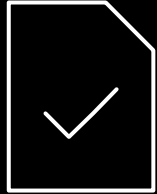
Julien Simon
Global Evangelist, AI & Machine Learning
Amazon Web Services
@julsimon

Stéphane Cheikh
Director, Portfolio Evolution using Artificial Intelligence
SITA

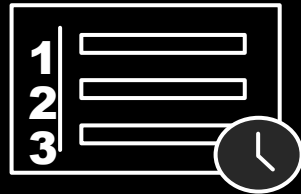
Machine learning cycle



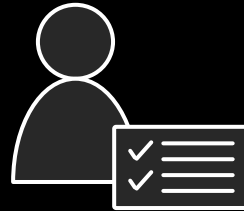
Amazon SageMaker



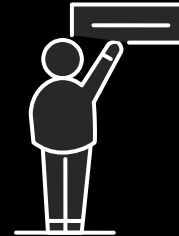
Collect and prepare
training data



Choose and
optimize your
ML algorithm



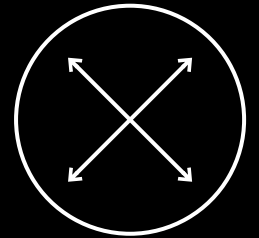
Set up and
manage
environments
for training



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



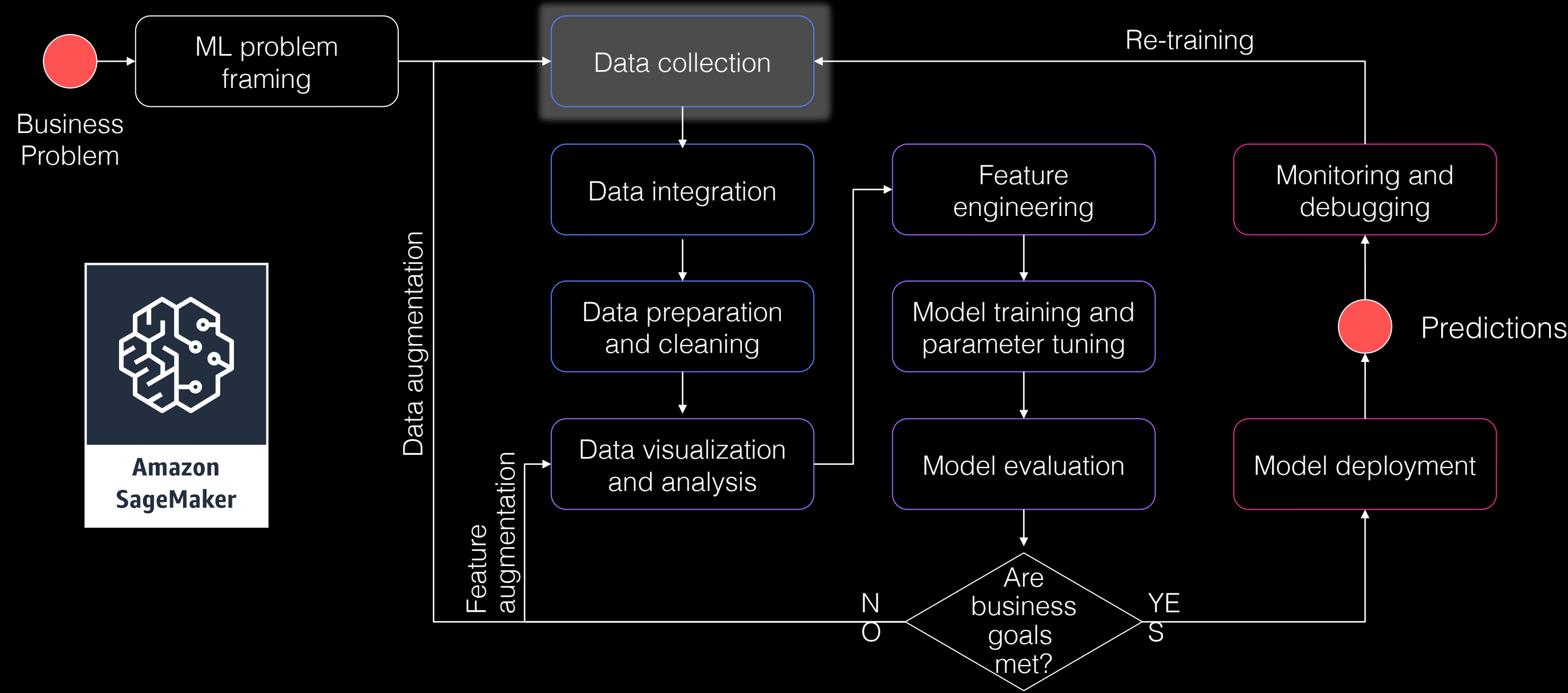
DOW JONES



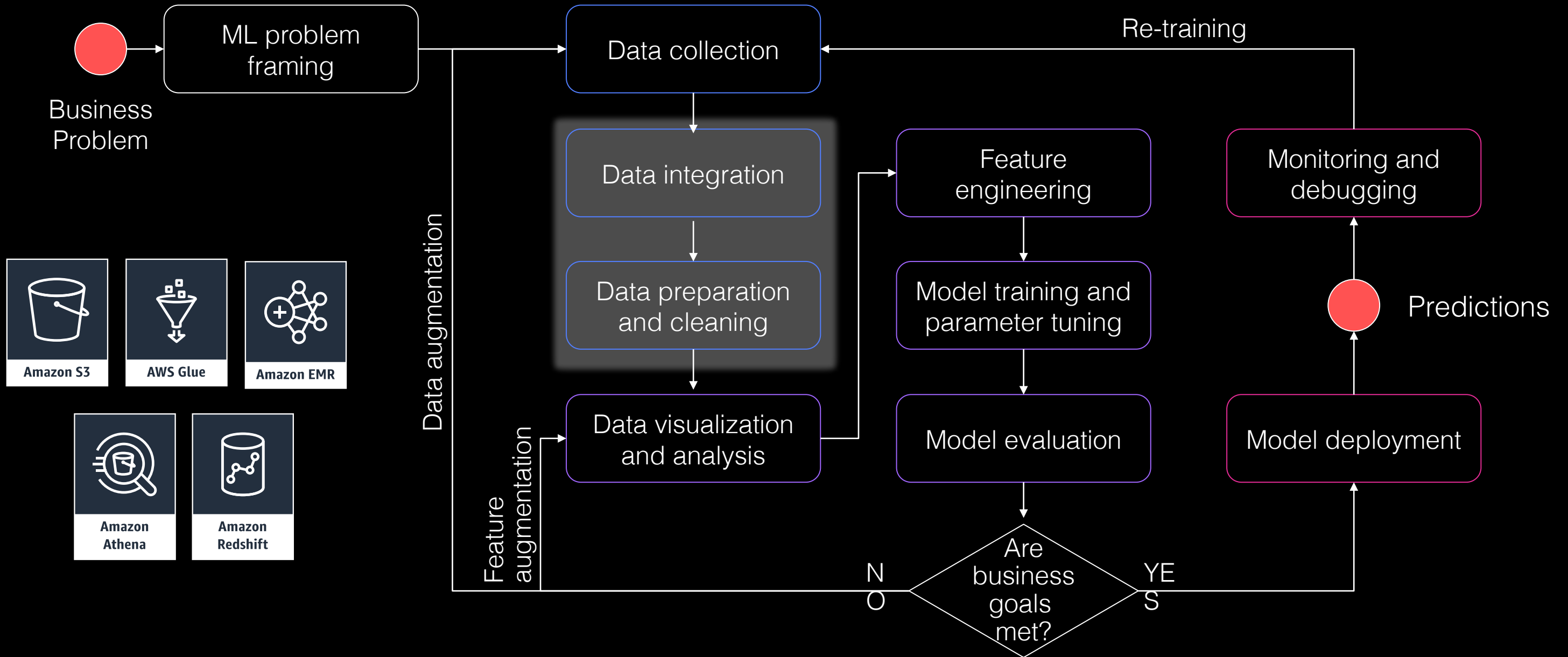
SONY



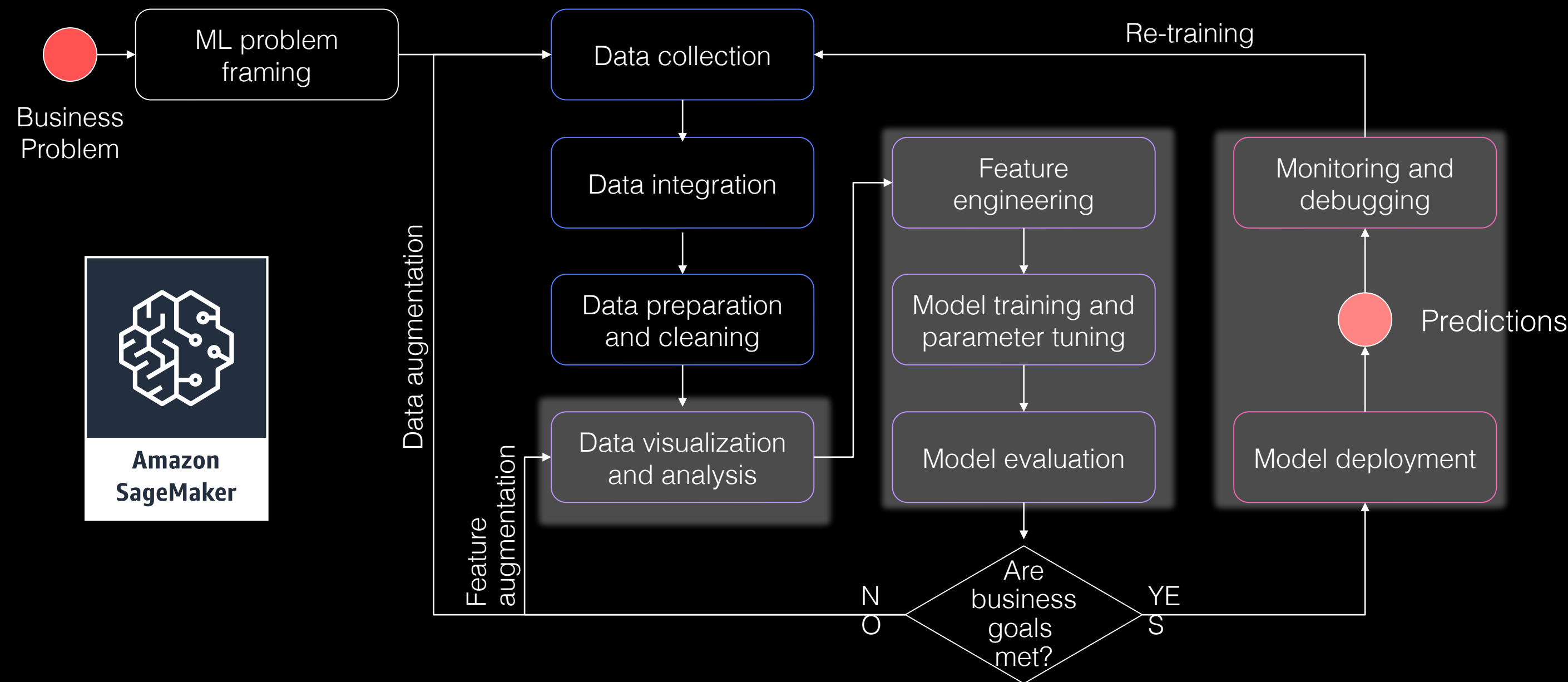
Build your dataset



Prepare your dataset for Machine Learning



Build, train and deploy models using SageMaker





WELCOME TO SITA

Flight Delay Prediction

Stephane Cheikh

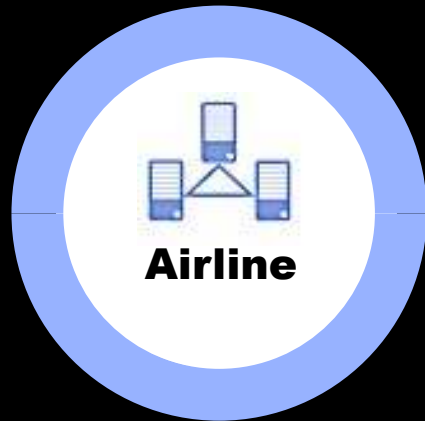
Director, Portfolio Evolution using Artificial Intelligence



SITA

SITA at a glance

AIR TRAVEL SOLUTIONS



Airline

Airline operations



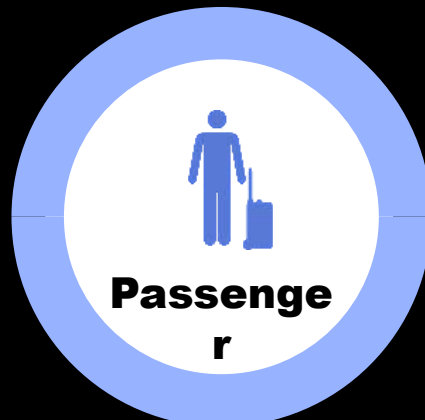
Airport

Airport operations,
Baggage processing,
Passenger processing



Government

Border
management



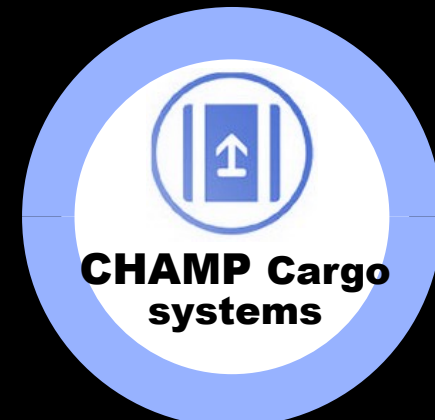
Passenger

SITA Passenger Service
System (SITA PSS)



SITAONAIR

Aircraft operations
(cockpit and cabin
services),
Connected aircraft



**CHAMP Cargo
systems**

Cargo management,
community integration,
eCargo

Disruption: Our Industry Issue

Airline Impact

- Delays
- Diversions & Cancellations
- Crew Time Limits
- Asset (Aircraft) Usage

Airport Impact

- Capacity
- Surges
- Out-of-Hours
- Extended Stays

Passenger Impact

- Baggage
- Accommodation
- Re-Booking
- Keeping Informed

Cargo Impact

- Perished Goods
- Storage Capacity
- Delivery Delays

76%

Global
OTP

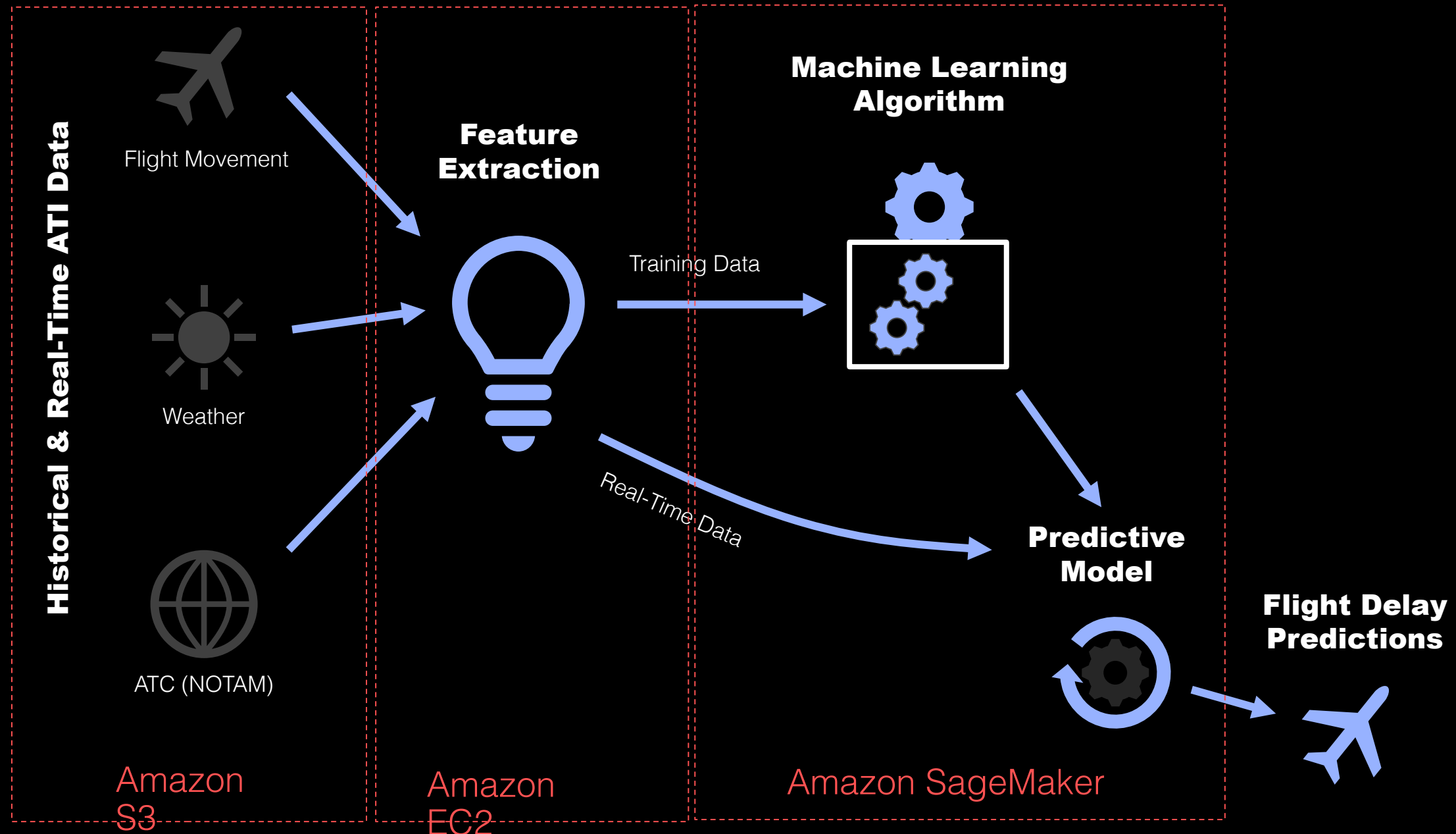
**51
mins**

Average
Flight Delay
Time

>\$25Bn

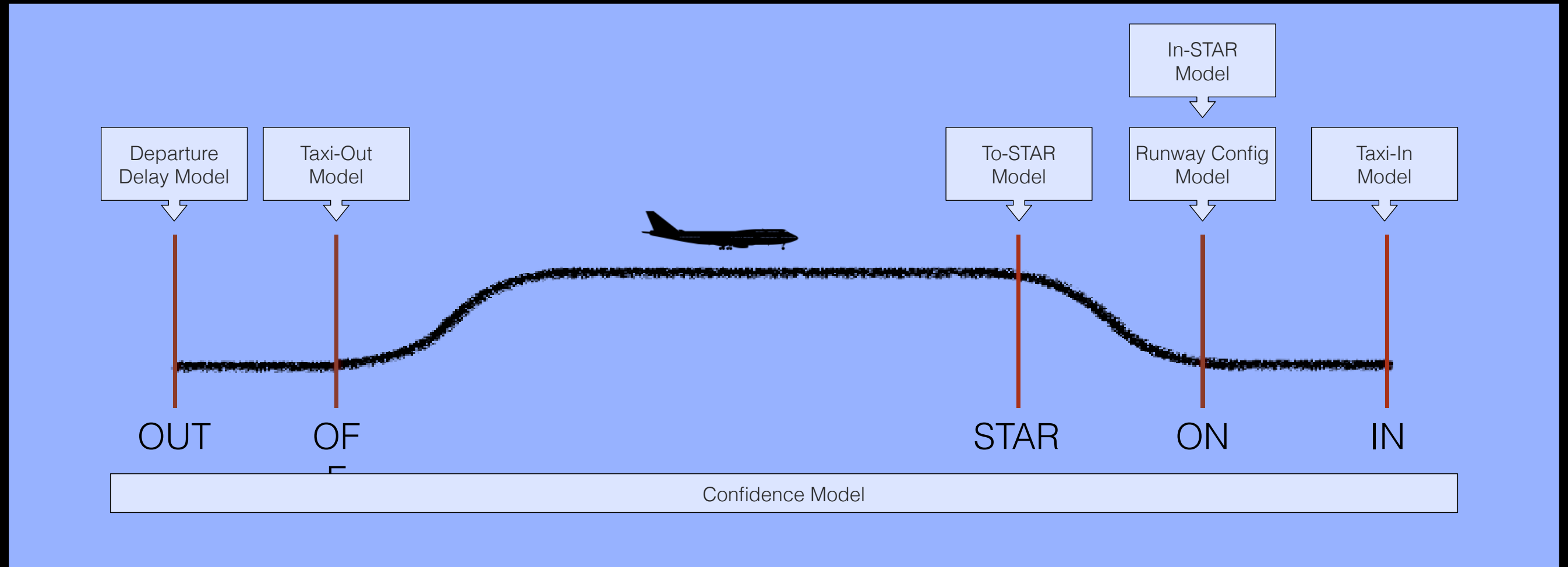
Global
Cost of
Delay

Applying AI to predict potential flight disruptions



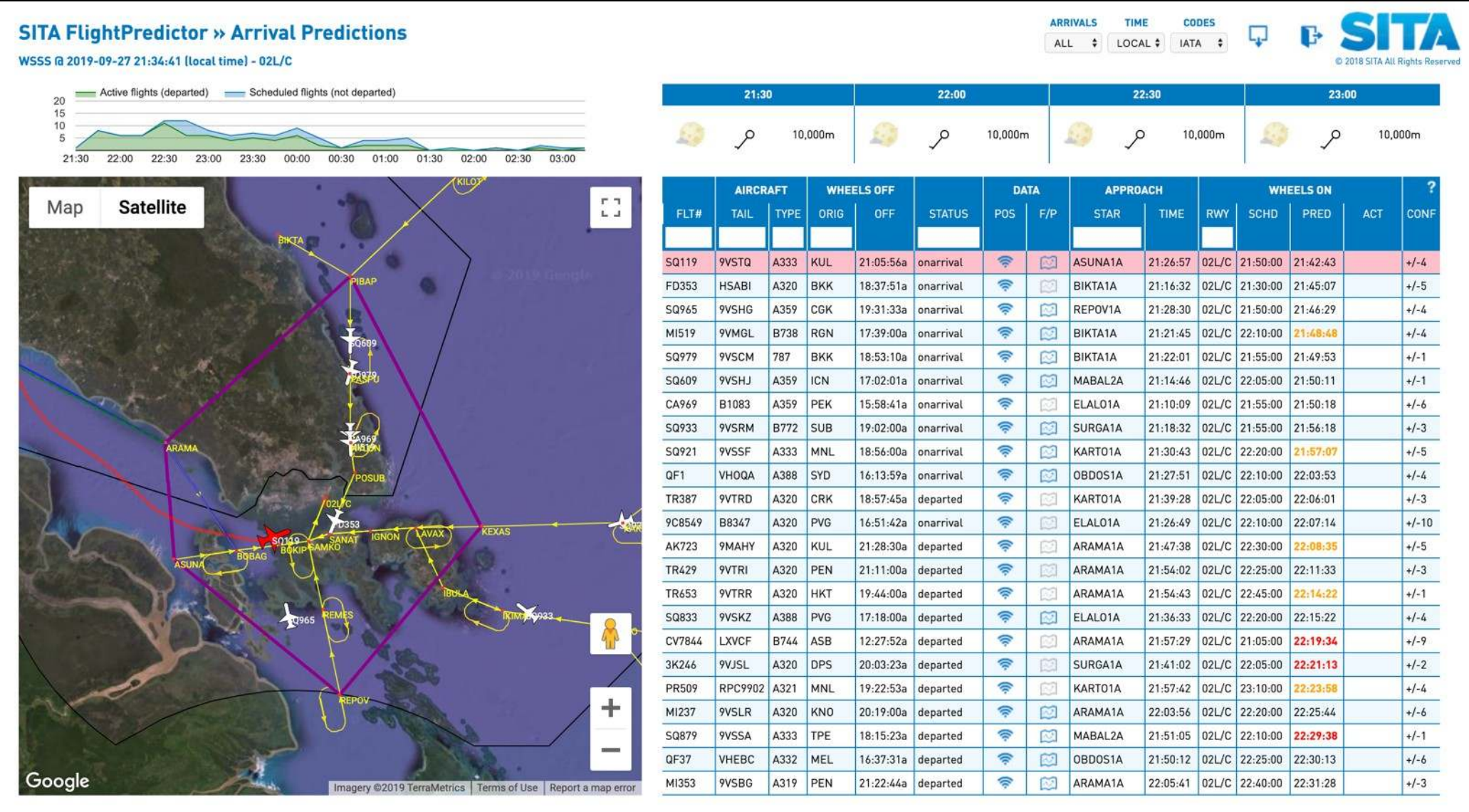
Seven models working together

All machine learning models re-usable and adaptable for different airports



SITA Flight Predictor

Long look-ahead predictions of flight schedule deviations for airline and airport users



What we achieved

Reliable and accurate predictions are possible up to 6 hours out

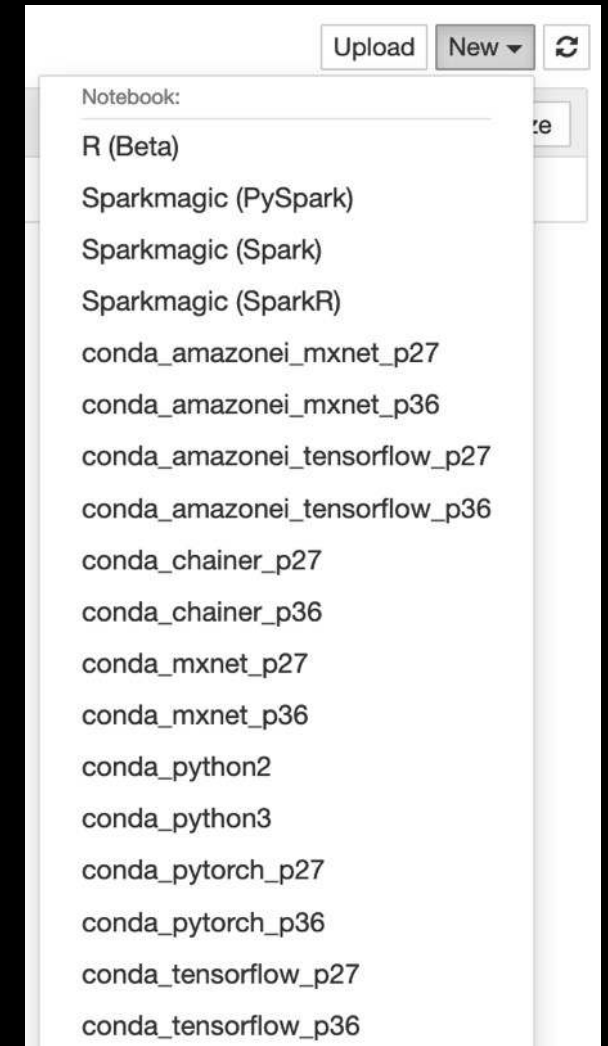
- “Fine tuning” will ensure consistent and reliable performance of all models
- Incorporating winds aloft data will improve predictions at 6 hours out
- New departure delay model will also improve accuracy of longer lead-time predictions

Prediction Time Interval	Prediction Accuracy
30 mins	±5 mins
1 - 2 hours	±10 mins
2 - 4 hours	±15 mins
6 hours	±15 mins

Building models

Notebook instances

- Fully managed EC2 instances, from *ml.t2.medium* to *ml.p3.16xlarge*
- 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**



Model options



Training code

**AWS Marketplace
for Machine
Learning:
250+ off-the-shelf
algos 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

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

Built-in algorithms

Built-in algorithms

Orange: supervised, yellow: unsupervised

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

Demo:

Sentence classification with BlazingText

https://github.com/aws-labs/amazon-sagemaker-examples/tree/master/introduction_to_amazon_algorithms/blazingtext_text_classification_dbpedia

Built-in frameworks

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**: Reuse **existing code** 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:

Fashion-MNIST 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.