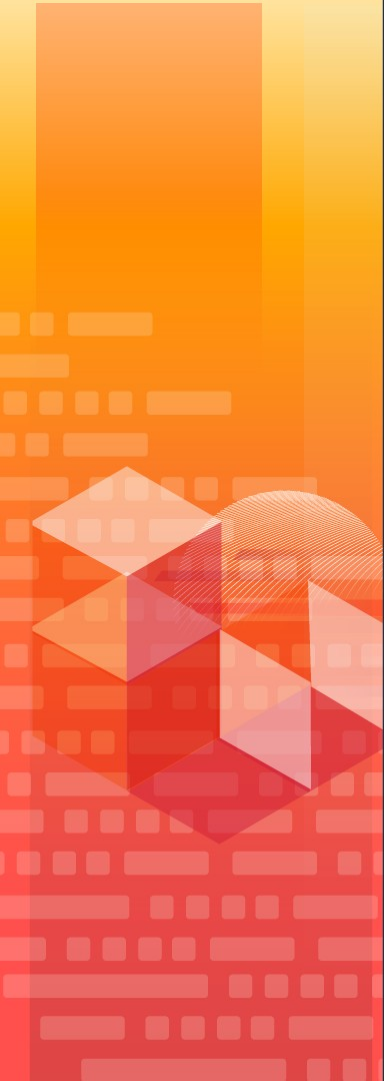# MLOps with serverless architectures

Julien Simon
Principal Technical Evangelist, AI and Machine Learning, AWS
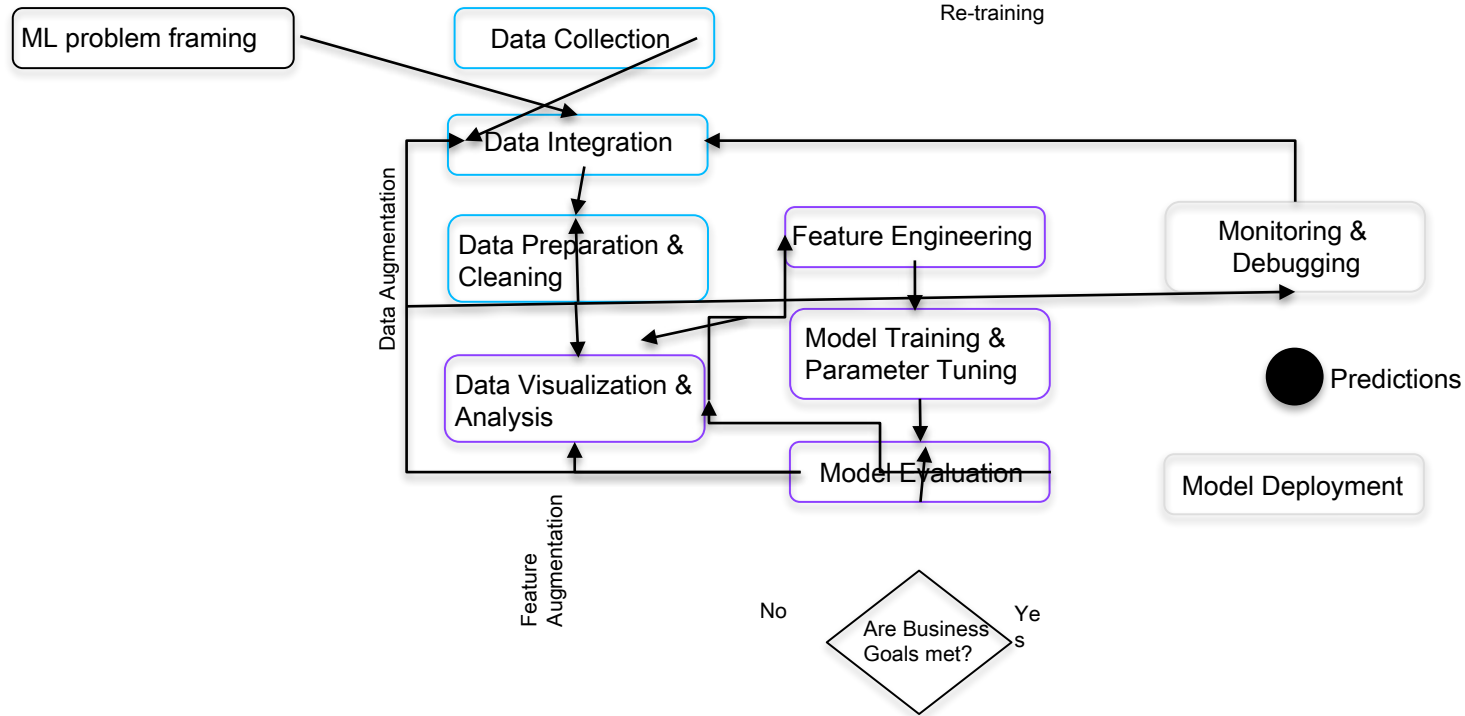
@julsimon

# Agenda

- But why?
- A quick recap on Amazon SageMaker
- A quick recap on serverless architectures
- Open Source tools: AWS Chalice, Serverless Framework
- Demos
- Resources

# Rationale

# The Machine Learning Process

Business Problem – ●

ML problem framing

Data Collection

Re-training

Data Integration

Data Augmentation

Data Preparation & Cleaning

Feature Engineering

Monitoring & Debugging

Data Visualization & Analysis

Model Training & Parameter Tuning

● Predictions

Model Evaluation

Model Deployment

Feature Augmentation

No     Are Business Goals met?     Yes

# Amazon SageMaker

| Build | Train | Deploy |
|-------|-------|--------|
| Pre-built notebooks for common problems | One-click training | One-click deployment |
| Built-in, high-performance algorithms | Hyperparameter optimization | Fully managed hosting with auto-scaling |

FREE TIER

# Ops needed?

Business Problem ●

ML problem framing

Data Collection

Data Integration

Data Preparation & Cleaning

Data Augmentation

Data Visualization & Analysis

Feature Augmentation

Feature Engineering

Model Training & Parameter Tuning

Model Evaluation

Re-training

Monitoring & Debugging

● Predictions

Model Deployment

Are Business Goals met?

No

Yes

# Serverless Architectures

# Serverless architecture

=

Fully-managed services

+

AWS Lambda



Amazon API Gateway

Amazon Kinesis Streams
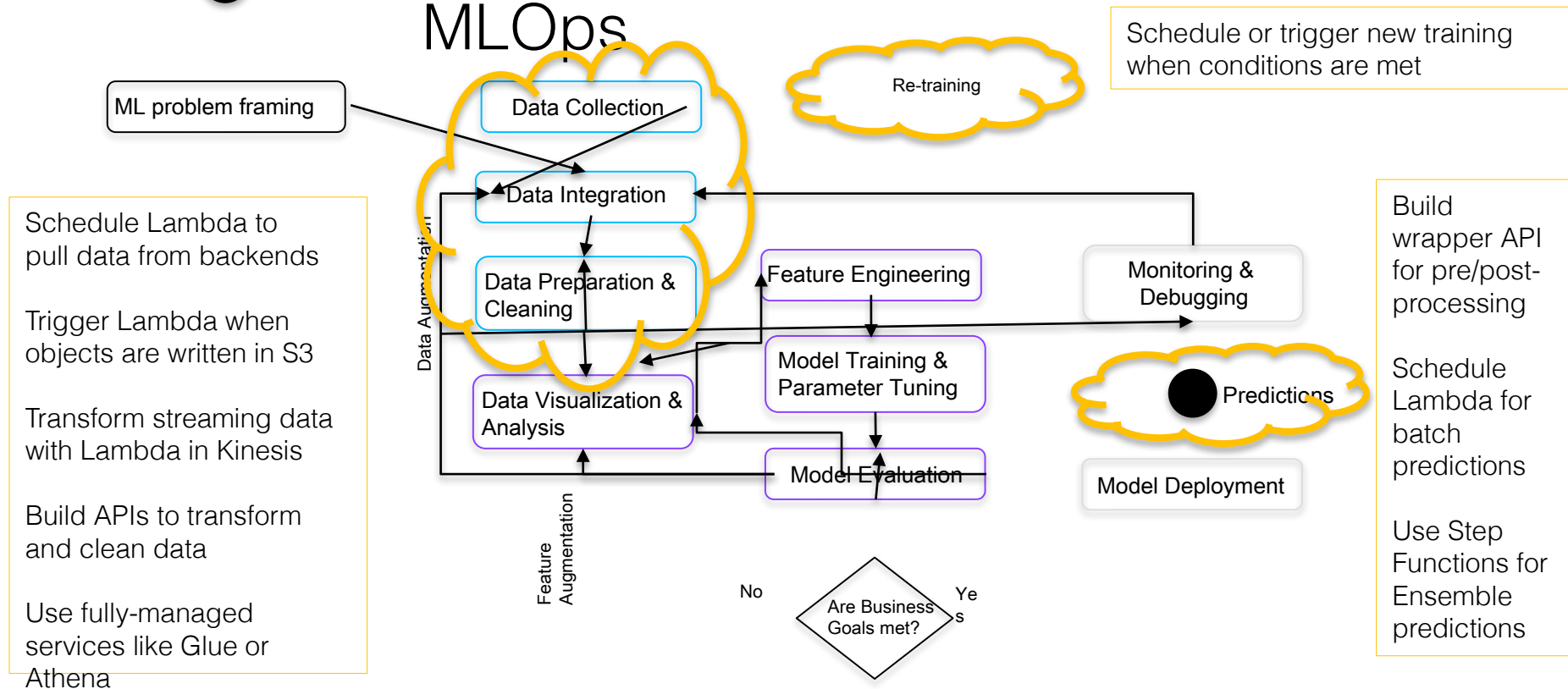
Amazon DynamoDB

Amazon S3

# AWS Lambda

- Announced at re:Invent 2014
- Deploy functions in Java, Python, Node.js ,C# and Go.
- Just code, without the infrastructure drama
- Built-in scalability and high availability
- Integrated with many AWS services
- Pay as you go
  - Combination of execution time (100ms slots) & memory used.
  - Starts at $0.20 per million requests.
  - Free tier available: first 1 million requests per month are free.
- Orchestration with AWS Step Functions.

# What can you build with serverless architectures?

- **Automate** your AWS infrastructure

- Build **event-driven** applications

- Build **APIs** together with Amazon API Gateway

# Ideas for serverless MLOps

**Business Problem** ●

ML problem framing

Data Collection

Re-training

Schedule or trigger new training when conditions are met

Data Integration

Data Augmentation

Data Preparation & Cleaning

Feature Engineering

Monitoring & Debugging

Build wrapper API for pre/post-processing

Data Visualization & Analysis

Model Training & Parameter Tuning

Predictions

Schedule Lambda for batch predictions

Feature Augmentation

Model Evaluation

Model Deployment

Use Step Functions for Ensemble predictions

No    Are Business Goals met?    Yes

Schedule Lambda to pull data from backends

Trigger Lambda when objects are written in S3

Transform streaming data with Lambda in Kinesis

Build APIs to transform and clean data

Use fully-managed services like Glue or Athena
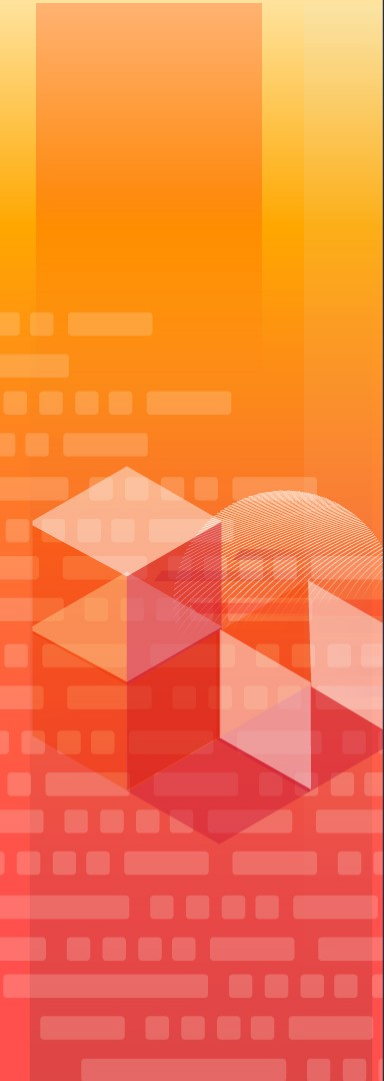
Open Source Frameworks

# The Serverless framework



- Announced at re:Invent 2015

- Auto-deploys and runs Lambda functions, locally or remotely

- Auto-deploys Lambda event sources: API Gateway, S3, DynamoDB, etc.

- Creates all required infrastructure with CloudFormation

- Simple configuration in YML

http://github.com/serverless/serverless
https://serverless.com

# AWS Chalice

- Open Source project released in July 2016
- "Flask for serverless"
- Just add your Python code
  - Deploy web services with a single CLI call and zero config
  - The API is created automatically
  - The IAM policy is auto-generated (crowd goes wild)
- Run and test APIs on local port 8000 (similar to Flask)

# Demo #1: resizing images with AWS Chalice

https://medium.com/@julsimon/using-chalice-to-serve-sagemaker-predictions-a2015c02b033

# Demo #2: building a prediction wrapper with AWS Chalice

https://medium.com/@julsimon/using-chalice-to-serve-sagemaker-predictions-a2015c02b033

# Demo #3: retraining models with the Serverless Framework

https://medium.com/@julsimon/retraining-sagemaker-models-with-chalice-and-serverless-71a585ddbc7d

# Resources

https://ml.aws
https://aws.amazon.com/sagemaker
https://aws.amazon.com/lambda

http://github.com/serverless/serverless
https://github.com/awslabs/chalice

https://medium.com/@julsimon

# Thank you!

Julien Simon

Principal Technical Evangelist, AI and Machine Learning

@julsimon