# AWS | Using AWS to Transform Customer Data in MongoDB into AI-driven Personalization

**Dylan Tong**, Machine Learning Partner Solutions Architect

aws machine learning

**Many possibilities...**

Forecasting

Customer Churn Prediction

**Personalization**

.
.
.

11000011
100011.....

Machine Learning

# Deep Convolutional Neural Networks

# Deep Reinforcement Learning

# Reinvent the Customer Experience

**Amazon Personalize**

- Session based recommendations

- Predictive Customer Analytics

**Amazon SageMaker RL**

- Contextual Bandits: uplift conversion rates

**Conversational AI**

- Humanize your apps with life-like voices

mongoDB.

aws machine learning

Recommender Modernization

# Classic Recommenders: Item-Item Collaborative Filtering



**1. Selection**  **2. Recommend "similar Item"**  **Other Products**

Calculate "rating vector" for each product and calculate vector distance to measure similarity

Customer 1

Customer N

Degree of similarity by co-rating

aws machine learning

# Deep learning techniques have a direct impact on the bottom line



**Popularity**

↓

**Matrix factorization**

**+15.4%**
Engagement

↓

**Neural network**

**+7.4%**
Engagement



**Similarity**

↓

**Recurrent Neural Net + Bandit**

**+20%**
Click Through

https://www.slideshare.net/AmazonWebServices/add-realtime-personalization-and-recommendations-to-your-applications-aim395-aws-reinvent-2018

|  10

aws machine learning

# RNN: History and User Representation

Customers interaction history: clicks, ratings, purchases…



recommend →



Unfold

$$o \quad o_{t-1} \quad o_t \quad o_{t+1}$$

$$W \quad W \quad W \quad W$$

$$V \quad h \quad \cdots \quad V \quad h_{t-1} \quad V \quad h_t \quad V \quad h_{t+1} \quad V \quad \cdots$$

$$U \quad U \quad U \quad U$$

$$x \quad x_{t-1} \quad x_t \quad x_{t+1}$$

→ User Representation

aws

# HRNN: Modeling Sessions

**Insight:** Evolution of interests and disinterests predict future preferences…



A month later…    recommend

Interactions, ordering and timing all matter…



Session Representations

User Representation

aws machine learning

# THE AWS ML STACK

Broadest and deepest set of capabilities

## AI Services

| VISION | | | SPEECH | | LANGUAGE | | CHATBOTS | FORECASTING | RECOMMENDATIONS |
|---|---|---|---|---|---|---|---|---|---|
| REKOGNITION IMAGE | REKOGNITION VIDEO | TEXTRACT | POLLY | TRANSCRIBE | TRANSLATE | COMPREHEND & COMPREHEND MEDICAL | LEX | FORECAST | PERSONALIZE |

## ML Services

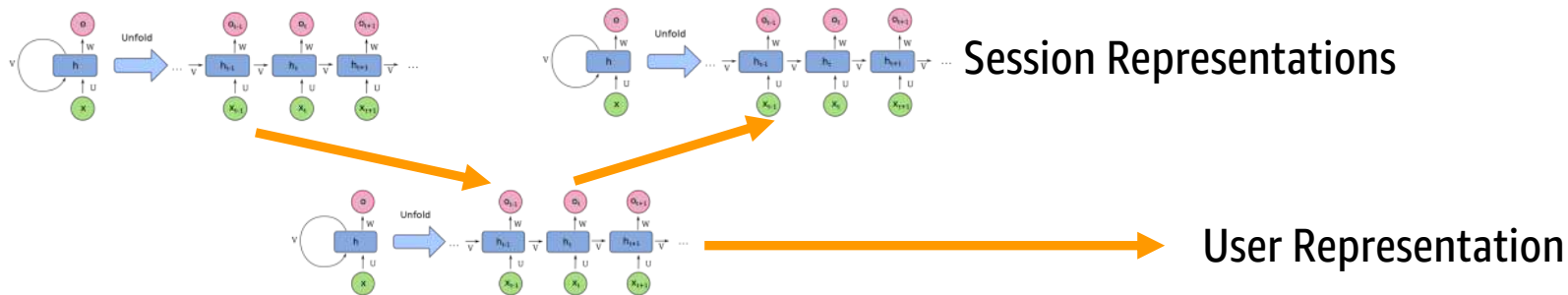| Amazon SageMaker | Ground Truth | Notebooks | Algorithms + Marketplace | Reinforcement Learning | Training | Optimization | Deployment | Hosting |
|---|---|---|---|---|---|---|---|---|

## ML Frameworks + Infrastructure

| FRAMEWORKS | INTERFACES | INFRASTRUCTURE | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| TensorFlow  mxnet  PYTORCH | GLUON  Keras | EC2 P3 & P3DN | EC2 G4 EC2 C5 | FPGAS | DL CONTAINERS & AMIs | ELASTIC CONTAINER SERVICE | ELASTIC KUBERNETES SERVICE | GREENGRASS | ELASTIC INFERENCE | INFERENTIA |

aws machine learning

# Amazon Personalize: AutoML

## Real-time Recommendations API



Events from App
Historical and
Online

Inventory

Demographics
(optional)

Amazon Personalize

INSPECT
DATA

IDENTIFY
FEATURES

SELECT
ALGORITHMS

SELECT
HYPERPARAMETERS

TRAIN
MODELS

OPTIMIZE
MODELS

HOST
MODELS

BUILD FEATURE
STORE

CREATE
REAL-TIME
CACHES

Customized
personalization &
recommendation
API

Fully managed by
Amazon Personalize

MongoDB

**GetRecommendations:**

```
{ "campaignArn": "string",
"itemId": "string",
"numResults": number,
"userId": "string" }
```

**GetRankedList:**

```
{ "campaignArn": "string",
"inputList": [ "string" ],
"userId": "string" }
```

aws machine learning

# Cold starts and Online Learning

Application



AWS SDK:
Event Recorder

```
personalize_events.put_events(
trackingId = 'tracking_id',
userId= 'USER_ID',
sessionId = 'session_id',
eventList = [{ 'sentAt': TIMESTAMP,
'eventType': 'EVENT_TYPE',
'properties': "{\"itemId\": \"ITEM_ID\"}" }])
```

MongoDB

Events from App
Historical and
Online

Inventory

Demographics
(optional)

Amazon Personalize

INSPECT
DATA

IDENTIFY
FEATURES

SELECT
ALGORITHMS

SELECT
HYPERPARAMETERS

TRAIN
MODELS

OPTIMIZE
MODELS

HOST
MODELS

BUILD FEATURE
STORE

CREATE
REAL-TIME
CACHES

Customized
personalization &
recommendation
API

Fully managed by
Amazon Personalize

aws machine
learning

# Predictive Customer Insights

- Predictive Customer-level Marketing

- Reverse Recommendations: query the users most likely to be interested in product(s).



MongoDB Atlas Data Lake Architecture

Batch Scores

MongoDB → Amazon S3 →

**Amazon Personalize** ✕

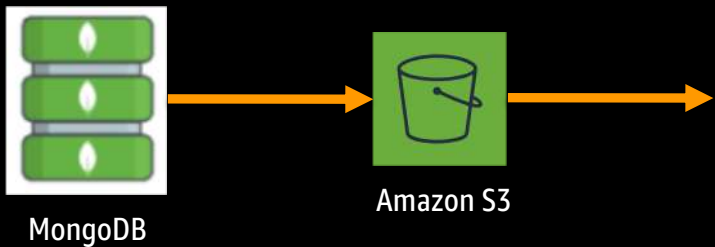Dataset groups

▼ **amazon-reviews**

Dashboard

Datasets

Event trackers

Solutions and recipes

Campaigns

## Dataset import job details

Dataset name

amazon-reviews

Schema name

amzn-reviews-interactions-v1

Schema

```
4      "namespace": "com.amazonaws.personalize.schema",
5      "fields": [
6        {
7          "name": "USER_ID",
8          "type": "string"
9        },
10       {
11         "name": "ITEM_ID",
12         "type": "string"
13       },
14       {
15         "name": "EVENT_TYPE",
16         "type": "string"
17       },
18       {
19         "name": "EVENT_VALUE",
20         "type": "string"
21       },
22       {
23         "name": "TIMESTAMP",
24         "type": "long"
25       }
```

Feedback   English (US)   © 2008 - 2019, Amazon Web Services, Inc.

aws machine learning

# Multivariate Optimization

# Contextual Bandits

After only a single week of online optimization, we saw a 21% conversion increase compared to the median layout...

# Reinforcement Learning (RL)

**State T**

**State T+1**

**Environment**

**Rewards**

**Policy Learner**

**Actions**

**Agent**

aws machine learning

# Multi-arm Bandit

Maximize expected outcome without knowledge of the true distribution.

Policy
Learner

Rewards

Actions

**Environment**
**77%**

**15%**

**25%**

# Contextual Bandits (CB)

**1. Context** (State)

- Demographics
- Device
- Promo Reference
- Geography

Agent

Policy Model

F**(context)** → **action**

**2. Actions**
(Select a layout)

aws machine learning

# Contextual Bandits: Multivariate Testing

**Arms = Layout Variations**



**Agent**

**1. Context** (State)

- Demographics
- Device
- Promo Reference
- Geography

Policy Model

F(**context**) → **action**

**2. Actions**
(Select a layout)

# Contextual Bandits and RL

**Application (Environment)**

3. Prescribe Layout

4. Rewards

Agent

**2. Actions**
(Select a layout)

**1. Context** (State)

Policy
Model

F**(context)** → **action**

# THE AWS ML STACK

Broadest and deepest set of capabilities

## AI Services

| VISION | | | SPEECH | | LANGUAGE | | CHATBOTS | FORECASTING | RECOMMENDATIONS |
|---|---|---|---|---|---|---|---|---|---|
| REKOGNITION IMAGE | REKOGNITION VIDEO | TEXTRACT | POLLY | TRANSCRIBE | TRANSLATE | COMPREHEND & COMPREHEND MEDICAL | LEX | FORECAST | PERSONALIZE |

## ML Services

| Amazon SageMaker | Ground Truth | Notebooks | Algorithms + Marketplace | Reinforcement Learning | Training | Optimization | Deployment | Hosting |
|---|---|---|---|---|---|---|---|---|

## ML Frameworks + Infrastructure

| FRAMEWORKS | INTERFACES | INFRASTRUCTURE | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| TensorFlow mxnet PYTORCH | GLUON Keras | EC2 P3 & P3DN | EC2 G4 EC2 C5 | FPGAS | DL CONTAINERS & AMIs | ELASTIC CONTAINER SERVICE | ELASTIC KUBERNETES SERVICE | GREENGRASS | ELASTIC INFERENCE | INFERENTIA |

aws machine learning

# Training Initial Model: Warm Starts (...*if data exists*)

**1. Experience data is prep and made available in the data lake.**

**Experience Data:**

**Source: web and application logs:**
- **Context features (state):** eg. device, geo, promo referrer...etc.
- **Action:** One of N layout variations
- **Action Probability:** chance that action is prescribed given the context for unbiasing the data.
- **Reward/Cost:** Selected value for a positive outcome. For instance, +1 for a click.

MongoDB Atlas

Amazon S3

MongoDB DataLake

aws machine learning

# Amazon SageMaker Training: "BYOS" Approach

**2. "Bring your own script":**

Algorithm

Environment

OR

Train

MongoDB Atlas

Amazon S3

MongoDB DataLake

Amazon SageMaker

# Bring Your Own Script for Vowpal Wabbit

**2. "Bring your own script":**



Train

Amazon S3

MongoDB DataLake

Amazon SageMaker

Vowpal Wabbit Contextual Bandits

Usage: ./vw -d train.dat **--cb_explore 10** --epsilon 0.1

**Amazon SageMaker Examples:**
VW Python Scripts (CLI wrapper)

aws machine learning

# Amazon SageMaker Training

3. **Launch Training Jobs:** SageMaker provisions a cluster and runs the training job—only pay for what you use.

```
estimator = RLEstimator(entry_point="train-vw.py",
                source_dir='src',
                dependencies=["common/sagemaker_rl"],
                image_name=custom_image_name,
                role=role,
                train_instance_type=instance_type,
                train_instance_count=1,
                output_path=s3_output_path,
                base_job_name=job_name_prefix,
                hyperparameters = {...}
        )

estimator.fit(...)
```

# Amazon SageMaker Hosting

## 4. Deploy the model for real-time inference

### I. Register model:

```
sagemaker_model = sagemaker.model.Model(
        image=self.image,
        role=self.resource_manager.iam_role_arn,
        name=model_id,
        model_data=model_record["s3_model_output_path"],
        sagemaker_session=self.sagemaker_session,
        env=environ_vars)
```

### II. Deploy endpoint:

```
sagemaker_model.deploy(
        initial_instance_count=hosting_instance_count,
        instance_type=hosting_instance_type,
        endpoint_name=self.experiment_id)
```

Export VW
model artifacts

Deploy
endpoint

SageMaker
Training

Amazon S3

SageMaker
Hosting

# Multivariate Testing in Production

**Application**

**MVT Service**

**1. Users' Context:**
{  Device: …
Geo: …
Promo: …}

**Policy Model**

**2. Action:**
Use Layout variant N

Real-time Endpoint
(Managed by
Amazon SageMaker)

MongoDB Atlas

# Exploration and Exploitation

# Exploration Policy

## Application



MongoDB Atlas

**1. Users' Context:**
{ Device: …
Geo: …
Promo: …}

**2. Action:**
Use Layout variant N
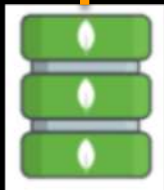
## MVT Service

### Exploration Policy:
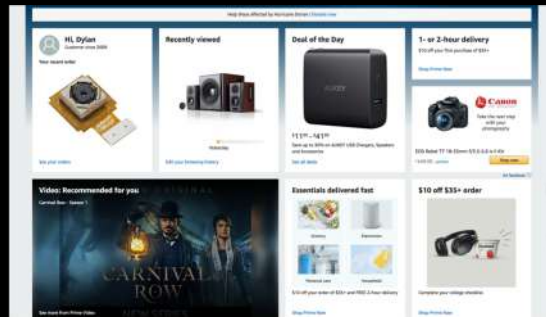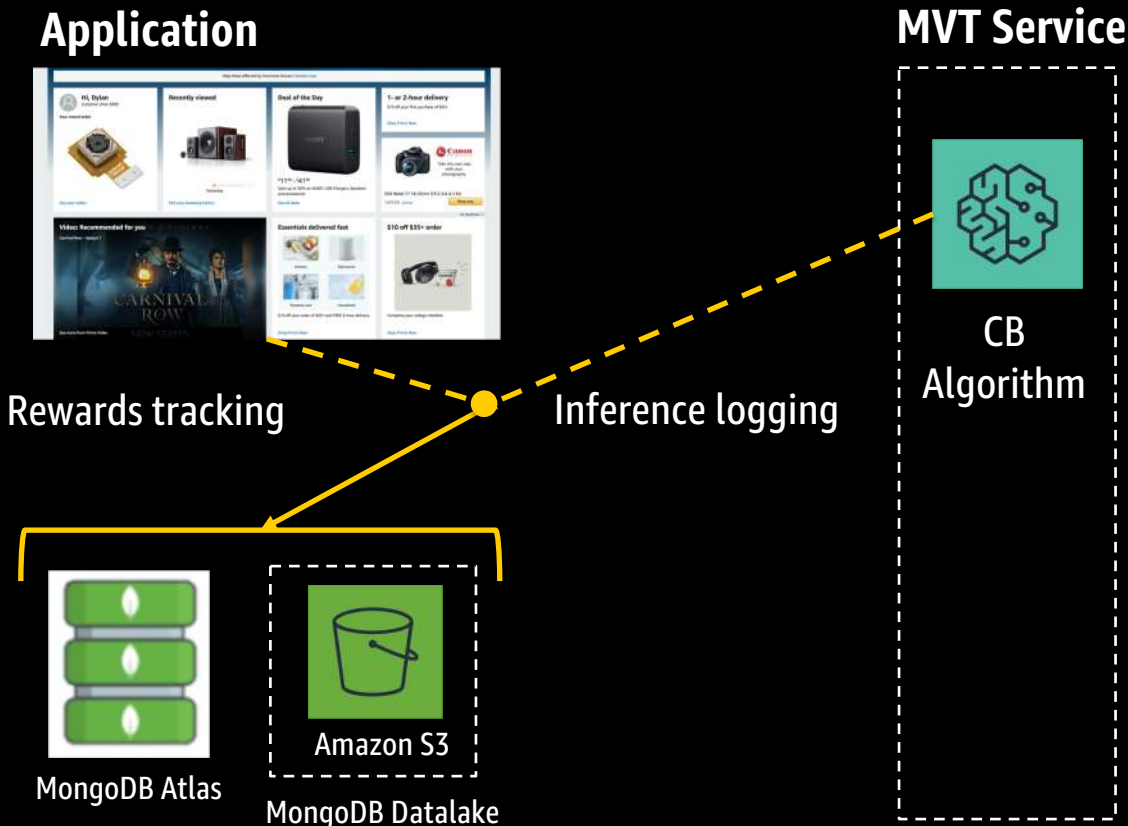
**Epsilon-Greedy:** Use action prescribed by trained policy model with probability (1-$e$), and one that is sampled uniformly at random with probability $e$.

Other policies: **UCB, Bagging, Online cover**…
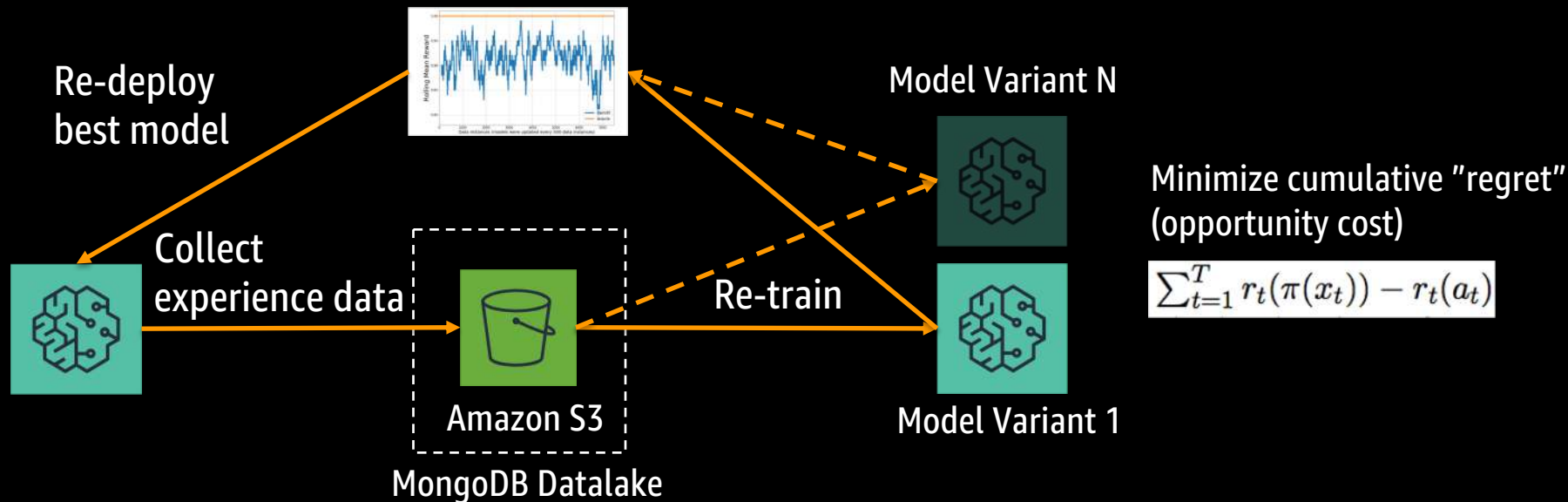
# Inference and Experience Capture

## 3. Capture Experiences:

I.  **Reward Tracking:** Event Id, Reward/Cost

II. **Inference Logging:** Event Id, Context, Action, Action Probability

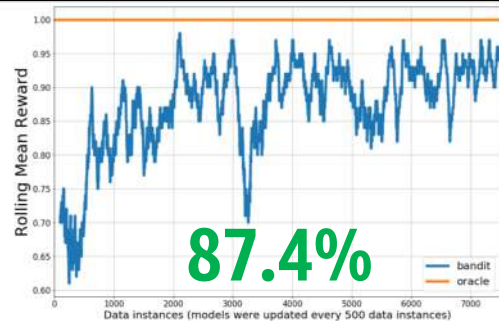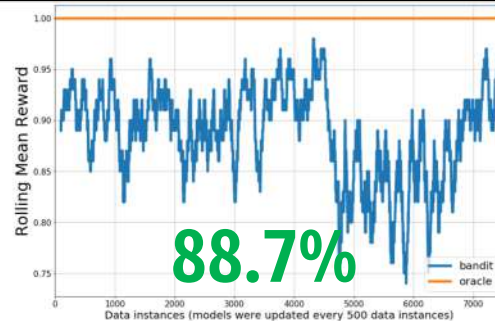III. **Associate** rewards with inference events to augment the training set (experience data).

**Application**

**MVT Service**



Rewards tracking

Inference logging

CB Algorithm

MongoDB Atlas

Amazon S3

MongoDB Datalake

aws machine learning

# Re-train, Evaluate and Re-deploy

Offline Evaluation (Replay): Compare new and old models with varying policies



Re-deploy
best model
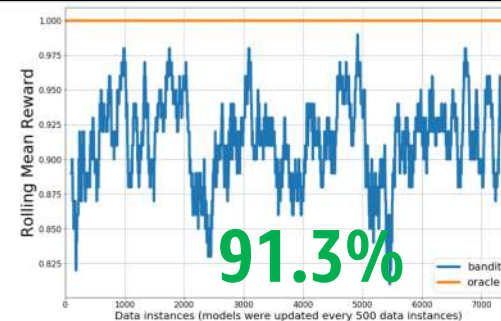
Model Variant N

Collect
experience data

Amazon S3

MongoDB Datalake

Re-train

Model Variant 1

Minimize cumulative "regret"
(opportunity cost)

$$\sum_{t=1}^{T} r_t(\pi(x_t)) - r_t(a_t)$$

aws machine
learning

# Converges Towards an Optimal Policy

# Humanize and Personalize Conversational AI

# THE AWS ML STACK

Broadest and deepest set of capabilities

## AI Services

| VISION | | | SPEECH | | LANGUAGE | | CHATBOTS | FORECASTING | RECOMMENDATIONS |
|---|---|---|---|---|---|---|---|---|---|
| REKOGNITION IMAGE | REKOGNITION VIDEO | TEXTRACT | POLLY | TRANSCRIBE | TRANSLATE | COMPREHEND & COMPREHEND MEDICAL | LEX | FORECAST | PERSONALIZE |

## ML Services

**Amazon SageMaker**

| Ground Truth | Notebooks | Algorithms + Marketplace | Reinforcement Learning | Training | Optimization | Deployment | Hosting |
|---|---|---|---|---|---|---|---|

## ML Frameworks + Infrastructure

| FRAMEWORKS | INTERFACES | INFRASTRUCTURE |
|---|---|---|
| TensorFlow  mxnet  PYTORCH | GLUON  K Keras | EC2 P3 & P3DN  EC2 G4 EC2 C5  FPGAS  DL CONTAINERS & AMIs  ELASTIC CONTAINER SERVICE  ELASTIC KUBERNETES SERVICE  GREENGRASS  ELASTIC INFERENCE  INFERENTIA |

aws machine learning

# Amazon Polly: Humanize Your Apps using Neural TTS



Sentence to synthesize.

'sɛntəns tə 'sɪnθə saɪz.

**Concatenative TTS**

'sɛnt    sɛntəns tə    'sɪnθ    ə saɪz.

Standard

US English Joanna voice

"President Donald Trump said on March 13 his administration was ordering the grounding of all Max 8 and 9 models, hours after Canada said it was grounding the planes after analyzing new satellite tracking data."

**Neural TTS**

DNN

Newscaster NTTS

aws machine learning

# Amazon Polly: Personalize Your Voices

**Justin**
English (US)
Male, Child

**Brian**
English (UK)
Male, Adult

Amazon Polly
Text-to-Speech
Lexicons
S3 synthesis tasks

## Text-to-Speech

Listen, customize, and download speech. Integrate when you're ready.

Type or paste your text in the window, choose your language and region, choose a voice, choose Listen to speech, and then integrate it into your applications and services.

With up to 3000 characters you can listen, download, or save immediately. For up to 100,000 characters, your task must be saved to an S3 bucket.

| Plain text | SSML | ❓ |

Welcome, to MongoDB London!

27 characters used

Show default text    Clear text

**Engine** ℹ️
- ○ Standard
- ● Neural

**Language and Region**

English, British ▾

**Voice**
- ○ Amy, Female
- ○ Emma, Female
- ● Brian, Male

▶ Listen to speech

⬇ Download MP3

Sample rate: 24000Hz
**Change file format**

Synthesize to S3

**Change S3 task settings**

aws machine learning

# Voice Modification

<speak>

This is Brian without any voice modifications.

<amazon:effect vocal-tract-length="+15%"> Imagine now that I got bigger... </amazon:effect>

<amazon:effect vocal-tract-length="+25%"> Suppose that I got even bigger still... </amazon:effect>

Now let's go back and hear the effect when I go in the opposite direction.

<amazon:effect vocal-tract-length="-15%"> Can you tell that I'm getting smaller? </amazon:effect>

<amazon:effect vocal-tract-length="-25%"> Now I'm even smaller than before. </amazon:effect>

</speak>

aws machine learning
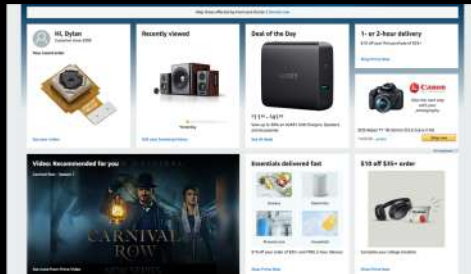
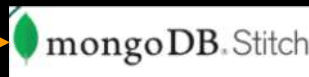# Deploy as a microservice

Your choice of compute…

Application



Public
API endpoints

Amazon
API Gateway

Backend service logic

mongoDB. Stitch

AWS Lambda

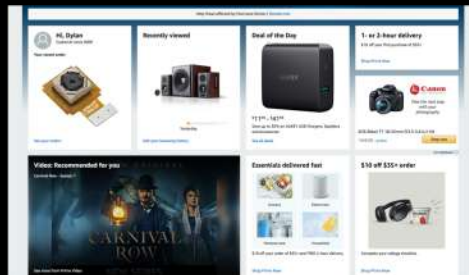Amazon Elastic
Kubernetes Service

· · ·

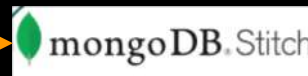aws machine learning

# Know your performance requirements

```
Params =
{ "Engine": "string",
"LanguageCode": "string",
"LexiconNames": [ "string" ],
 "OutputFormat": "string",
"SampleRate": "string",
"SpeechMarkTypes": [ "string" ],
"Text": "string",
"TextType": "string",
"VoiceId": "string" }
```
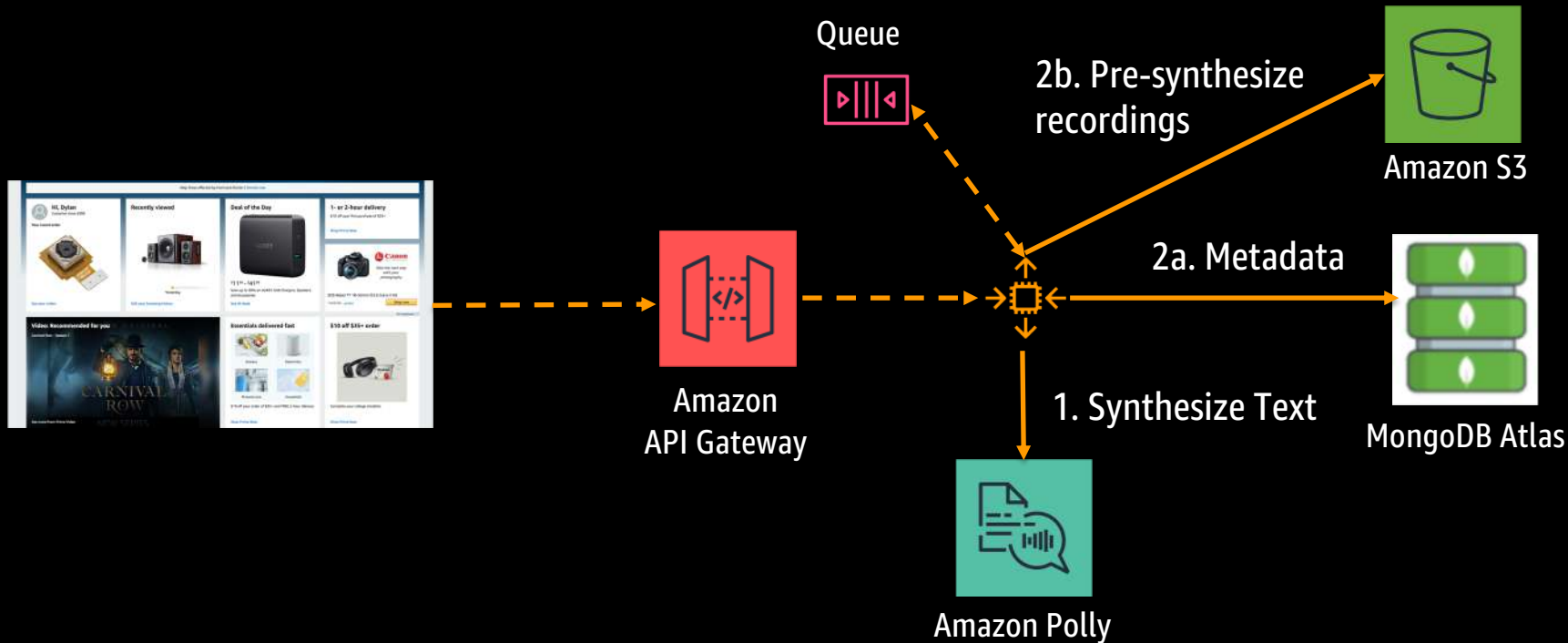
Polly.synthesizeSpeech
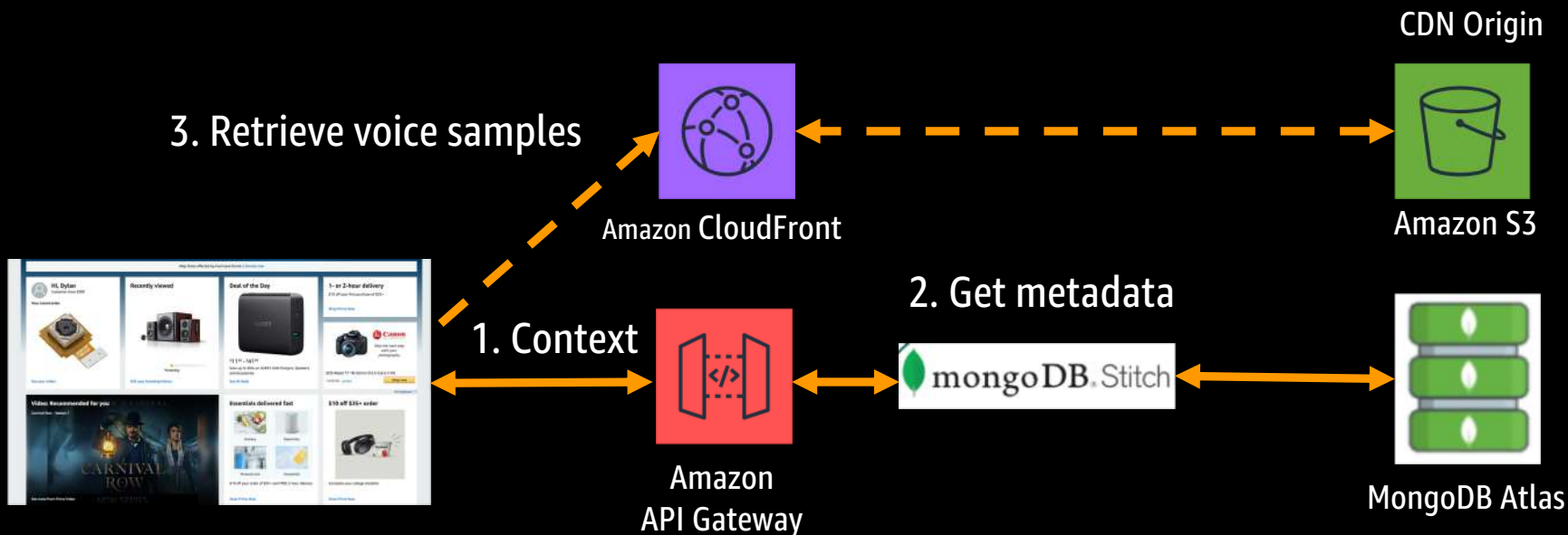      (params, (err, data) => {…}

**Amazon
API Gateway**

mongoDB.Stitch

**80 TPS
100 TPS (burst)**

**Amazon Polly**

aws machine learning

# Build in Caching and Pre-processing



Queue

2b. Pre-synthesize recordings

Amazon S3

2a. Metadata

MongoDB Atlas

Amazon
API Gateway

1. Synthesize Text

Amazon Polly

aws machine learning

# Cache and stream from the edge



3. Retrieve voice samples

CDN Origin

Amazon CloudFront

Amazon S3

1. Context

Amazon
API Gateway

2. Get metadata

mongoDB.Stitch

MongoDB Atlas

aws machine learning

# Branding Voice: chat bots with personality



## Github: Lex Chatbot Example