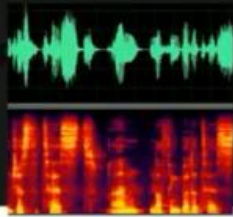Image recognition is a field of deep learning that uses neural networks to recognize the subject and traits for a given image. In Japan, Cookpad uses Amazon ECS to run an image recognition platform on clusters of GPU-enabled EC2 instances. In this session, hear from Cookpad about the challenges they faced building and scaling this advanced, user-friendly service to ensure high-availability and low-latency for tens of millions of users.

# Machine learning

Significantly improve many applications on multiple domains

"deep learning" trend in the past 10 years

aws



# Machine learning at Amazon
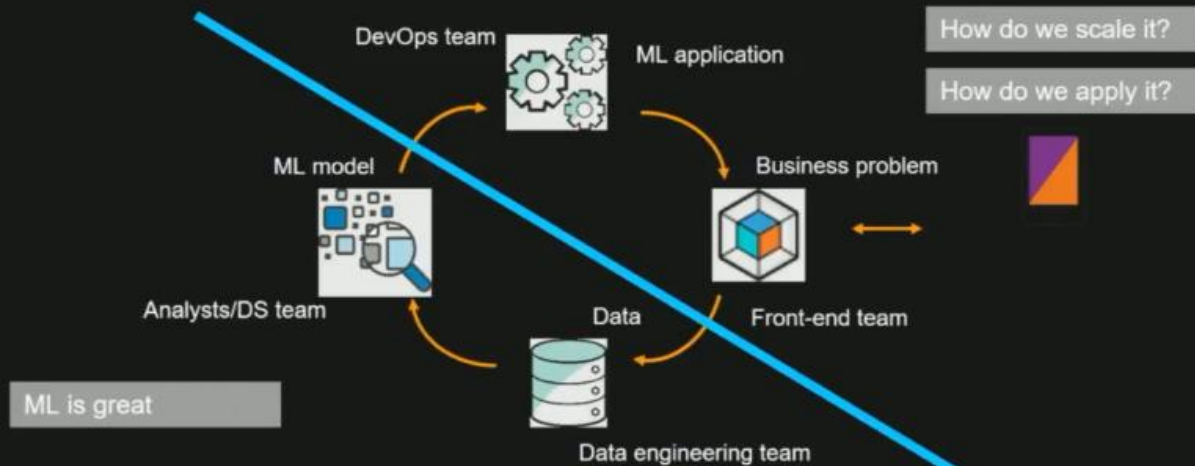
JUST ASK

NO LINES. NO CHECKOUT (NO, SERIOUSLY.)

# The circle of machine learning (ML)

DevOps team

ML application

ML model

Business problem

How do we scale it?

How do we apply it?

Analysts/DS team

Data

Front-end team

ML is great

Data engineering team

# Solution for building AI apps with CICD

# AWS AI services

| Amazon Rekognition | Amazon Polly | Amazon Lex | AI Services |
|---|---|---|---|
| Amazon Machine Learning | Amazon EMR | Spark & Spark ML | AI Platforms |

| Apache MXNet | TensorFlow | Caffe | Torch | Theano | CNTK | Keras | AI Engines |
|---|---|---|---|---|---|---|---|

The other primitive we need to build our service is a container, we will need to deploy 1000s of containers on 1000s of nodes, we can use ECS for this.
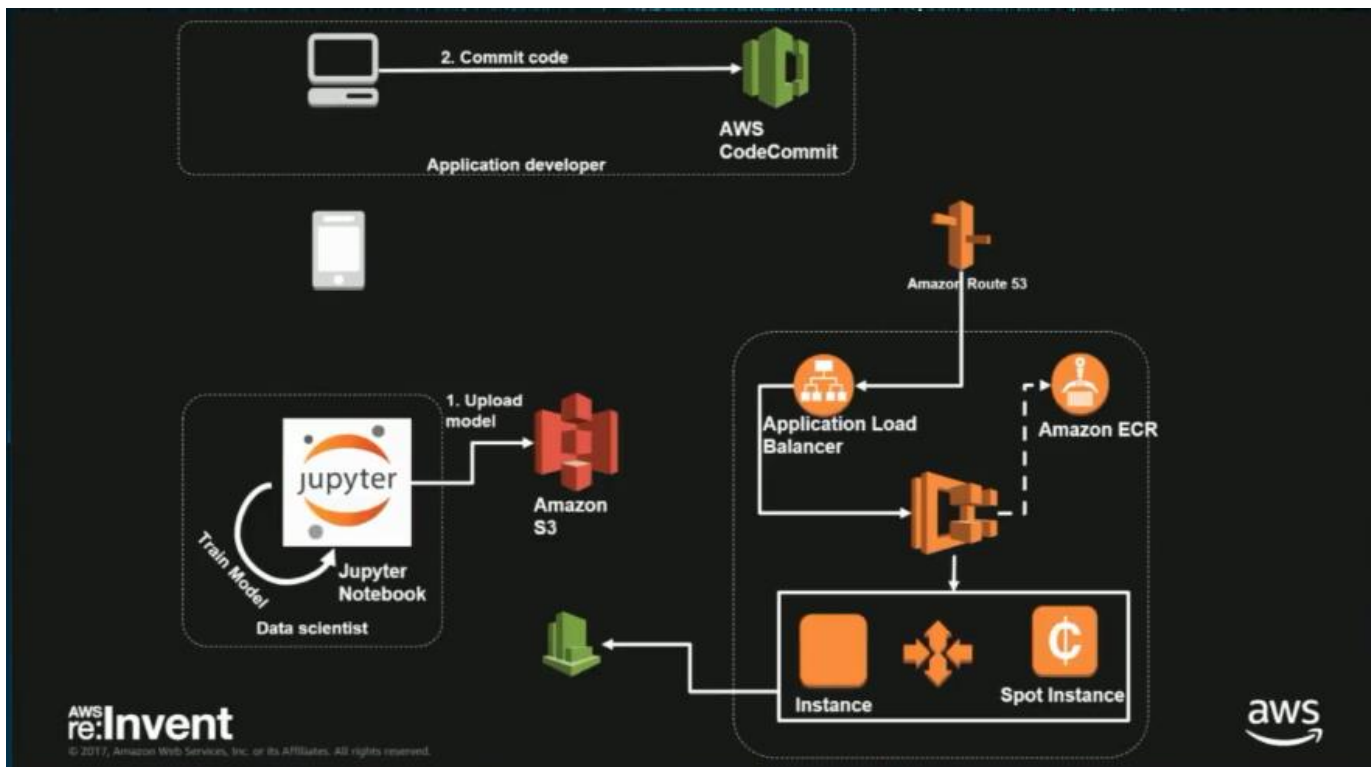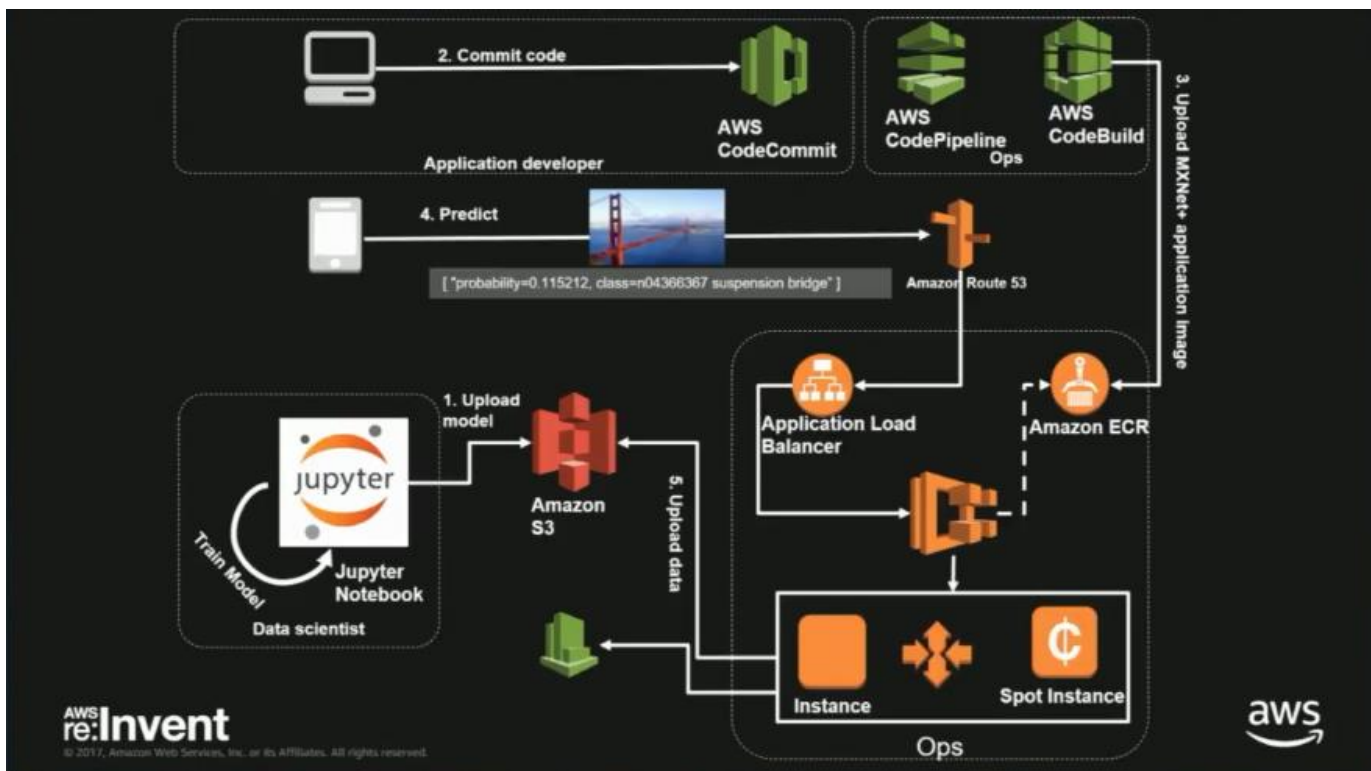
Amazon has published AMIs that you can just click to run in an EC2 instance and you can use for training your models and doing work as a data scientist. We upload our models into S3.



The application developer wants to write code, push it to a repo, build and deploy the app.

We also need a compute to run our workloads, you can download a CF template that can do this for you today

# Give it a spin



http://amzn.to/2zWpQij

aws

- Hokuto Hoshi (@kani_b)
- Head of Infrastructure, Cookpad Inc.
- hokuto@cookpad.com



aws ✓ CERTIFIED
🎁 Solutions Architect - Professional
🎁 DevOps Engineer - Professional

# About Cookpad

- "Make everyday cooking fun!" - Since **1998**
- https://cookpad.com/
- Largest online recipe sharing and search service in Japan

About **60M** Monthly users in Japan

Over **2.7M** user-authored recipes

# Cookpad is global

- https://cookpad.com/#{your_country_code}
  - us, uk, id, es, fr, br, ae, etc....

**67** countries

**21** languages

Offices in Japan, UK, Spain, Indonesia etc.

---

# Our infrastructure



**150+** developers

**9** SREs
**9** Machine learning engineers

All-in on AWS since **2011**

powered by **aws**

**~1,400** EC2 instances
**200+** ECS services

Over **2** regions

**15,000+** requests/sec

---

# Cooking Log（料理きろく）

**Our very first Deep Learning powered feature**

# Cooking log
# (料理きろく)

Collect "Food Photos" from
Camera Roll automatically

Powered by
Convolutional Neural Network

---



# DEMO

2017年11月

料理きろく 🔒

写真を読み込んでいます（残り36枚）

写真を読み込むのに、しばらく時間がかかる場合があります



アプリトップに戻る

ホーム　　タイムライン　　投稿する　　お知らせ　　MYキッチン

**140,000+**
Users

**12,000,000+**
"food" photos

# Our first deep learning feature

- There were several new challenges
    - Semi real-time image classification in production
    - Different workloads from the rest of web applications
- Especially we needed:
    - Scalable infrastructure for new workloads
    - Environment isolation for new challenge

# Scalable and asynchronous classification

- What we needed: Scalable infrastructure for massive photo uploading and semi real-time classification
  - Clients send tiny thumbnails after taking photos (difficult to predict traffic)
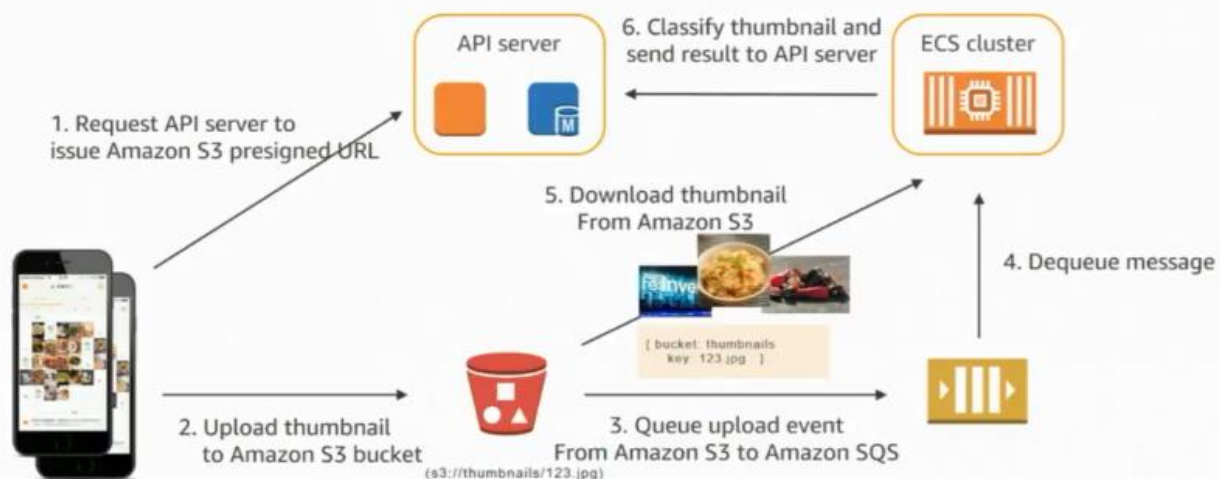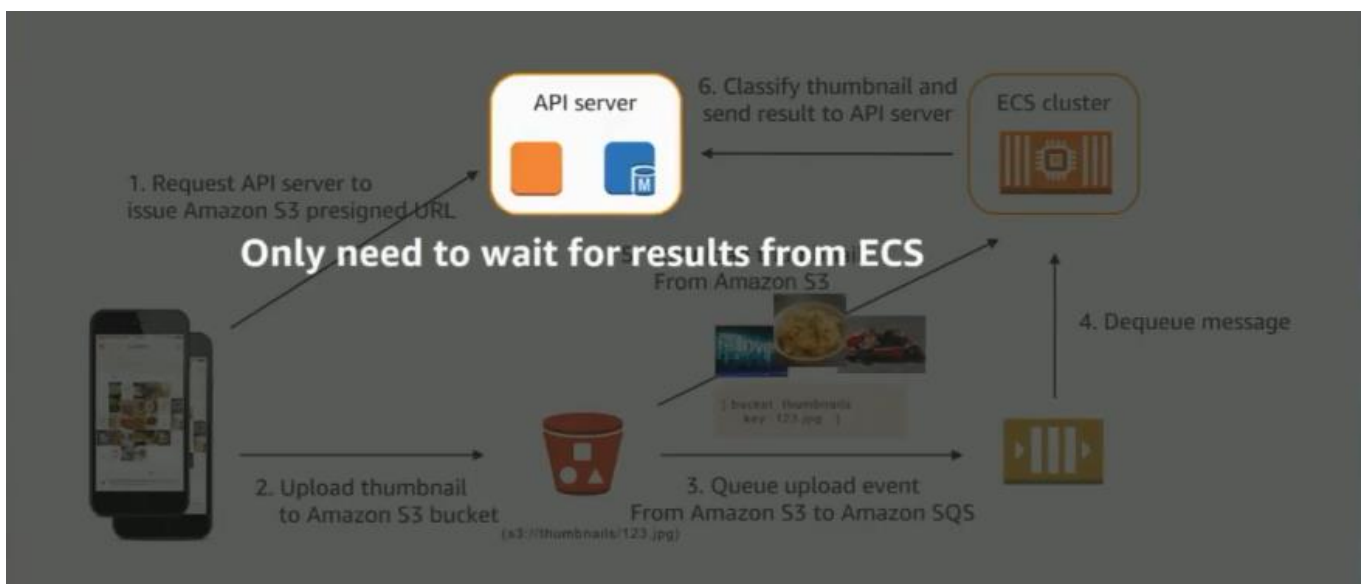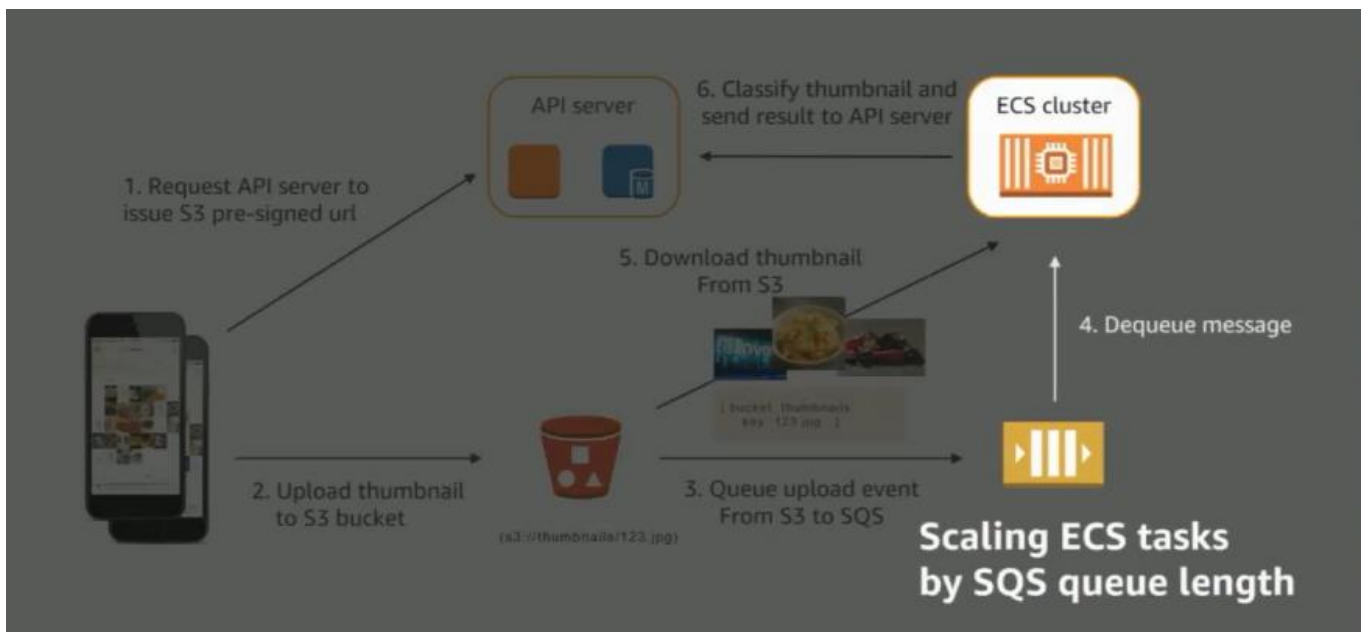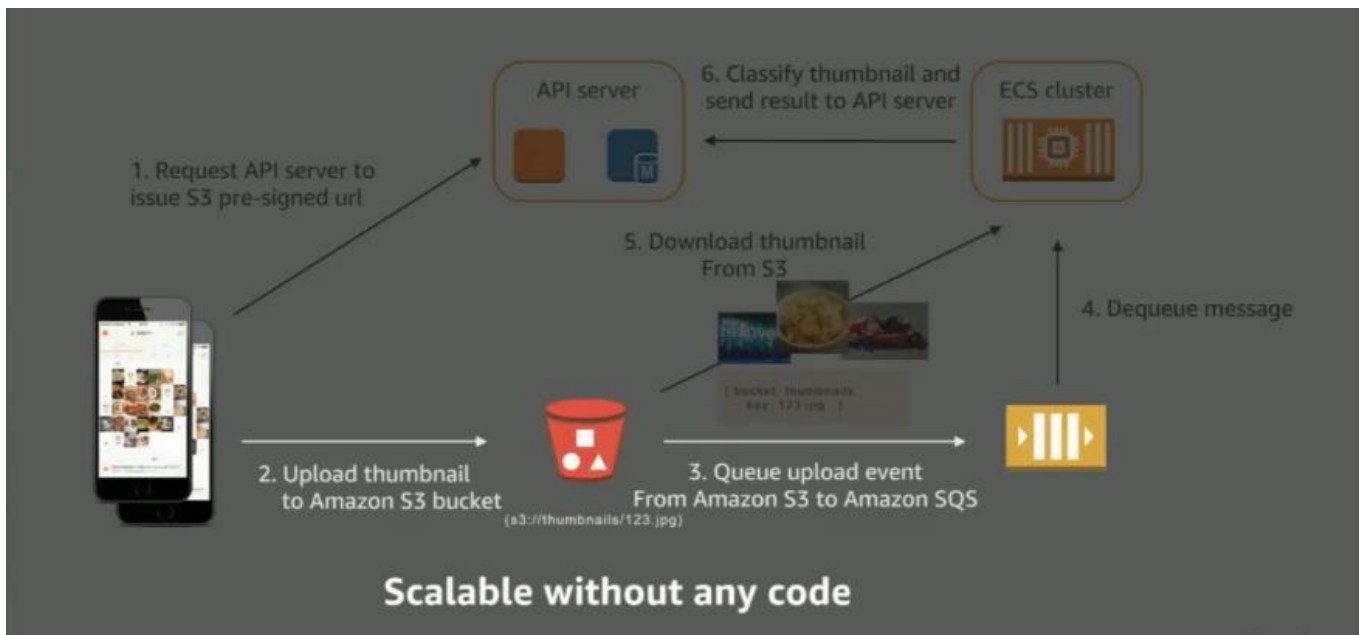  - Traffic spikes are coming sometimes (e.g. The TV show introduces our app)

- What we chose: Asynchronous architecture
  - Uploading and classification take time (~ several hundred ms)
    - Synchronous processing with API servers gives users a bad experience
  - Upload directly from clients to Amazon S3 using presigned URL
  - Enable Amazon S3 notification and Amazon SQS for queue of classific

# Environment isolation

- What we needed: Isolated environment in production
  - Different languages (Python for Machine Learning, Ruby for web application)
  - Different workloads (It's our first product that uses deep learning)
  - Different hardware (GPU)

- What we chose: Container environment
  - The container ensures runtimes are isolated
    - Language environment, GPU drivers, many configurations
  - Amazon ECS provides managed and scalable Docker environment
  - And we had already used containers on Amazon ECS!
  - We run all classification in Amazon ECS cluster on g2.xlarge

API server

6. Classify thumbnail and send result to API server

ECS cluster

1. Request API server to issue Amazon S3 presigned URL

5. Download thumbnail From Amazon S3

4. Dequeue message

{ bucket: thumbnails
  key: 123.jpg  }

2. Upload thumbnail to Amazon S3 bucket

(s3://thumbnails/123.jpg)

3. Queue upload event From Amazon S3 to Amazon SQS

Scalable without any code



Scaling ECS tasks by SQS queue length



Only need to wait for results from ECS
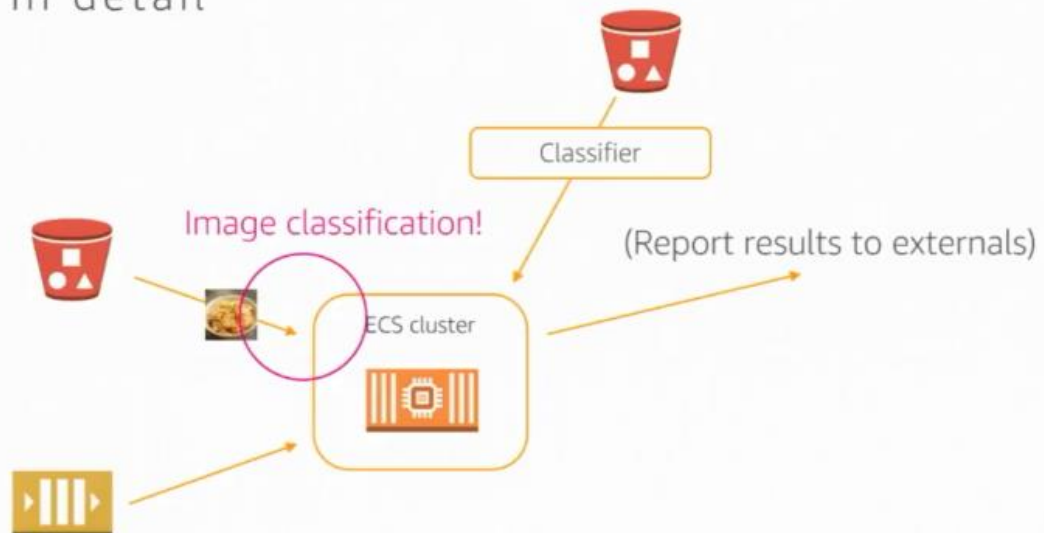
# Machine learning and infrastructure

Infrastructure to accelerate our machine learning projects.

- Yuichiro Someya (ayemos)
- Machine Learning Engineer @ Cookpad Inc.
    - # 2016(new grads) ~ Current

**aws** ✓ **CERTIFIED**
Solutions Architect - Associate
Developer - Associate

# The task in detail

Image classification!

Classifier

(Report results to externals)

ECS cluster

# Deploy the classifier on Amazon ECS

(Serialized) classifier

ECS cluster

Task definition?

- Fetch the classifier
- Dequeue from Amazon SQS
- Load the image from Amazon S3
- Report the results back

docker

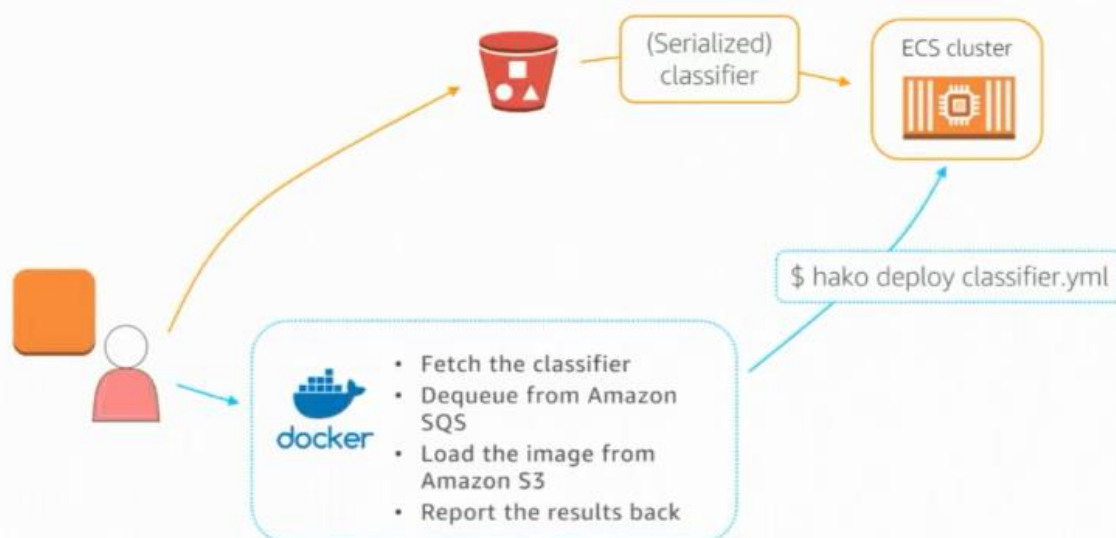# Democratize task definitions

- hako: https://github.com/eagletmt/hako/
  - Container deploy tool (Amazon ECS compatible)
  - Use yaml as definition file format

- Each developer writes app.yml and sends Pull request
- **200+** applications are at work

```yaml
scheduler:
  type: ecs
  region: ap-northeast-1
  cluster: hako-production-g2
  desired_count: 1
app:
  image: food-photo-classifier
  cpu: 128
  memory: 3072
  memory_reservation: 2048
  env:
    AWS_REGION: ap-northeast-1
      COOKPADNET_ENV: production
        ...
```
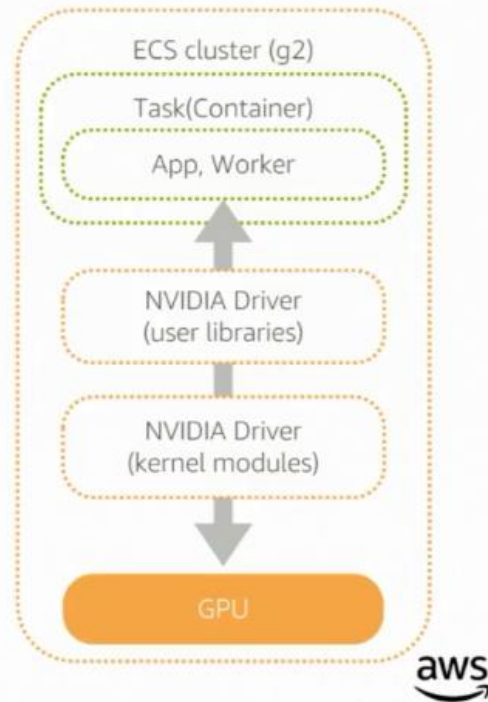
# Why we're using hako

- `hako` behaves as abstract operation layer over Amazon ECS (or other docker manager)
  - Higher-level operations: Deploy/Rollback/Stop/Remove
  - Manages *secret* environment variables
  - Pluggable pre/post development operations as `scripts`
    - Operations like DNS settings, Consul registrations, and so on
- We want each developer can deploy tasks on Amazon ECS individually.
  - `hako` handles Infra/SRE work around Amazon ECS.

# Deploy the model on Amazon ECS



- (Serialized) classifier
- ECS cluster
- $ hako deploy classifier.yml
- Fetch the classifier
- Dequeue from Amazon SQS
- Load the image from Amazon S3
- Report the results back

aws

# ECS and GPU

- Install drivers to clusters
  (ref: https://github.com/NVIDIA/nvidia-docker/wiki/NVIDIA-driver )
- CUDA to the container
  (CUDA ⇔ Driver version compatibility is relatively loose)

- GPU device files have to be visible and writable
  - **Privileged** flag
    (migrating to linux_parameters.devices option)
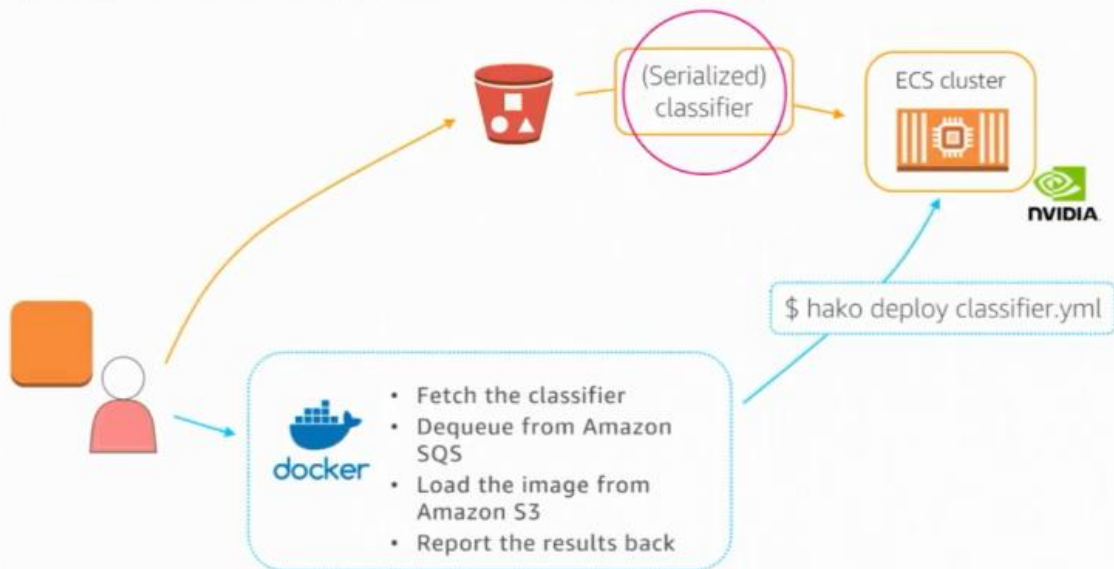
ECS cluster (g2)

Task(Container)

App, Worker

NVIDIA Driver
(user libraries)

NVIDIA Driver
(kernel modules)

GPU

# Deploy the model on Amazon ECS

(Serialized) classifier

ECS cluster

NVIDIA

$ hako deploy classifier.yml

docker
- Fetch the classifier
- Dequeue from Amazon SQS
- Load the image from Amazon S3
- Report the results back

Labeled image dataset

Food photos from Cookpad

train

Random photos from other datasets (Nonfood)

Food/Nonfood image classifier

Whole set

Food (100,000 photos~)

Nonfood (100,000 photos~)

?

$\overline{\{food\}} \neq \{non\text{-}food\}$

It's hard to obtain "Complement set"

Whole set

Food

Nonfood

?

{Food: 50%, Nonfood: 50%}

Whole set

Food

Nonfood

Plushies

Rebuild dataset
making use of posterior insights

Labeled image dataset

Food Photos from Cookpad

Random Photos from other datasets (Nonfood)

Food/Nonfood Image classifier

(Precision: 97.6%, Recall: 95.8%)

97.9% Accurate

aws





HashiCorp Packer

CUDA 8.0/cuDNN7

New!

CUDA 9.0/cuDNN7

`ssh workbench-001`

aws

## Cooking log

- Scalable food/non-food image classification
  - Asynchronous and isolated architecture
  - Containerized GPU workloads
- `hako` makes it easy to define and deploy applications

## Machine learning infrastructure

- Great environment makes our research fast and creative!
  - Managing multiple AMIs using Packer
  - Dedicated Instances
  - Operate instances via chat bot

# We're hiring!

https://info.cookpad.com/us

https://github.com/cookpad

cookpad

aws