




# Amazon SageMaker Ground Truth

第 6 回 Amazon SageMaker 事例祭り

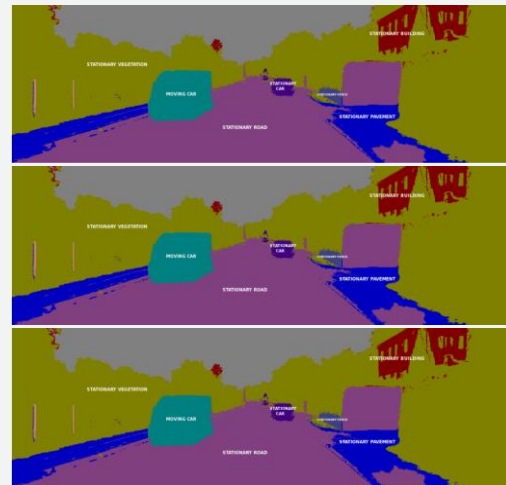
 #sagemaker\_fes

アマゾン ウェブ サービス ジャパン株式会社

# 独自のデータを利用した機械学習の流れ



大量の高品質な  
ラベル付けが重要



# 独自のデータを利用した機械学習の流れ



**アノテーション（データへのラベル付け）にはコスト・時間がかかる**

- 進捗管理・作業割り振り
- 効率の良いラベリングツールの作成
- 作業を割り当てるワーカーの募集
- これらを用意した上で数万個のデータへのラベル付け…



**これらの課題を解決するのが Amazon SageMaker Ground Truth**

# Amazon SageMaker Ground Truth

New!

## データにラベル (Ground Truth) を付与するアノテーション作業の支援サービス

- アノテーションの一般的なワークフローをサポート
- 4種類の組み込みラベリングツールを提供
- アノテーション作業を行うワーカーとの連携・管理機能を提供
- 大規模データセットに対しては自動ラベリング機能で最大70%のコスト削減



迅速なラベル付け



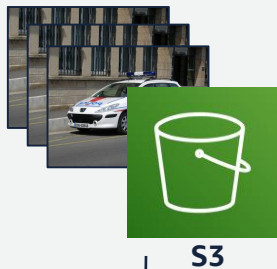
ワーカーとの連携が容易



高精度

# Ground Truth 利用のワークフロー

## 1. アノテーション対象のデータをアップロード

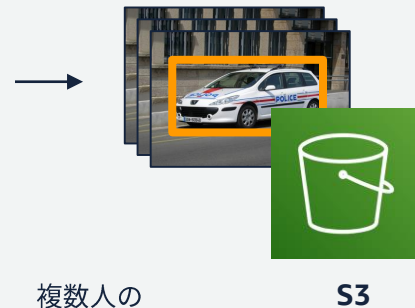


## 3. タスクはワーカーに自動で割り振られる

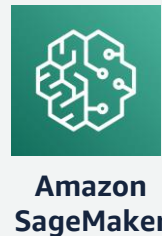
## 4. ラベリングツールでワーカーがアノテーションを行う



## 5. アノテーション結果が S3に格納される



## 6. 学習・推論に利用



## 2. ラベリングジョブの作成



複数人の  
結果をマージ

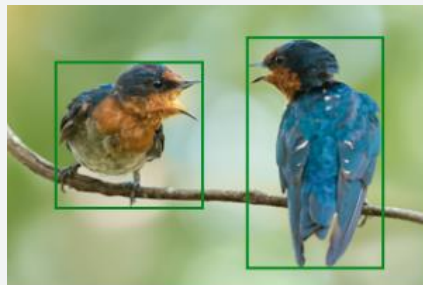
S3

ビルトインアルゴリズムを  
そのまま適応可能

# 組み込みラベリングツールの利用も 独自実装も可能



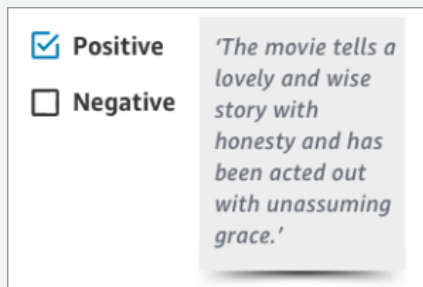
画像分類



物体検出



セマンティック  
セグメンテーション



文章分類



カスタム

# 組み込みのラベリングツールの使用例

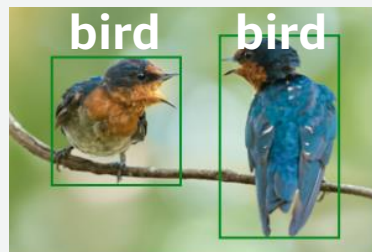
## 画像分類（画像から1つのラベルを推定）

- 撮影した画像から商品を推定（ラベルは100種類の商品）
- 顔画像から感情を推定（ラベルは喜怒哀楽+無表情の5種類）
- 独自の基準での不適切画像判定（ラベルは適切/不適切の2種類）



## 物体検出（画像から物体のラベルと矩形領域を推定）

- 製品や道路の画像から損傷箇所を検出
- レントゲン写真からある病気の箇所を検出する



# 組み込みのラベリングツールの使用例

## セマンティックセグメンテーション（ピクセル単位のラベリング）

- 走行画像の分析（車や人、走路情報などの領域を抽出）



## 文章分類（文章から1つのラベルを推定）

- レビュー文章のポジティブ/ネガティブ判定
- 文章のカテゴリ分類（政治・スポーツ・芸能などのラベル）

※ 分類対象の文章には日本語をセット可能  
（日本語が理解できるワーカーが必要なことに注意）

<input checked="" type="checkbox"/> Positive	<i>'The movie tells a lovely and wise story with honesty and has been acted out with unassuming grace.'</i>
<input type="checkbox"/> Negative	



# 組み込みのラベリングツール

画像分類

# ワーカーは以下の3種類から選択可能



## パブリック

- クラウドソーシングサービスの Amazon Mechanical Turk を利用
- 非言語依存で機密性の低いタスク向き

## プライベート



- 友人や社員をワーカーとして登録出来る
- 機密性の高いタスク向き
- ワーカーの管理にCognitoを利用（SAMLでの連携も可）

## ベンダー



- SageMaker Ground Truthに[登録済みのアノテーション専門ベンダー](#)に依頼
- 現時点では日本のベンダーは登録されていない

# 自動ラベリング

データの一部をワーカーがラベル付けするだけで、  
残りのラベル付けが自動化され、時間とコストを大幅に削減



※ 5000データ以上の大規模データセットに対して利用可能なオプション機能

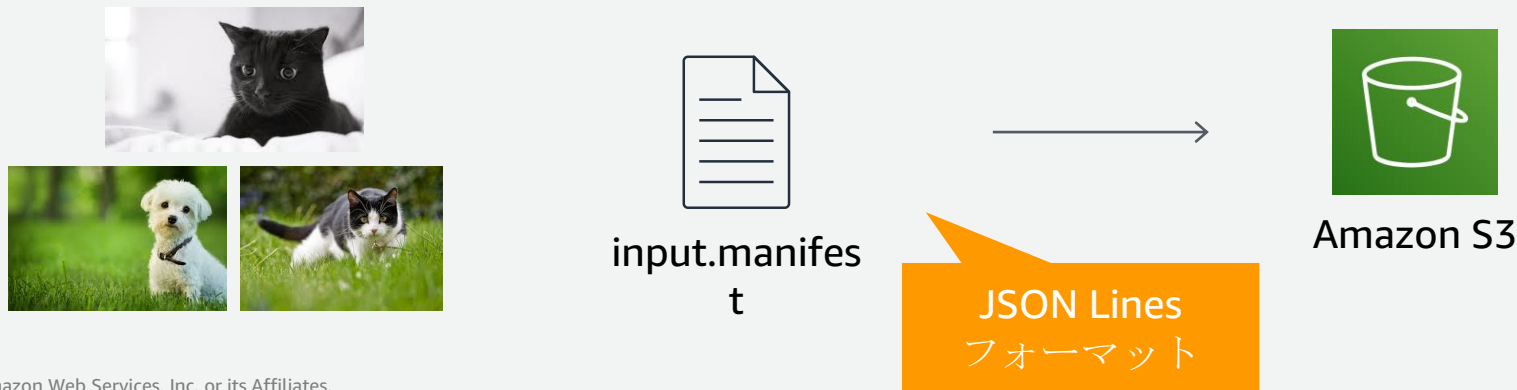
# ラベリングジョブの作り方

# ラベリングジョブ作成の4ステップ

1. データセットの準備
2. タスクの定義
3. ワーカーの選択
4. ラベリングツールの設定

# 1. データセットの準備

- データと結果を保存するための**S3**バケットを用意
- アノテーション対象の画像群を**S3**バケットに保存  
(テキスト分類の場合は **Text/CSV** ファイル)
- 画像群のパスを記述したマニフェストファイルを作成し、**S3**バケットに保存  
(マニフェストファイルは自動生成も可能)



# 1. データセットの準備 (マニフェストファイル)

input.manifest

```
{"source-ref": "s3://sagemaker-groundtruth-demo/SSDB00001.JPG"}  
{"source-ref": "s3://sagemaker-groundtruth-demo/SSDB00002.JPG"}  
{"source-ref": "s3://sagemaker-groundtruth-demo/SSDB00003.JPG"}  
{"source-ref": "s3://sagemaker-groundtruth-demo/SSDB00004.JPG"}  
{"source-ref": "s3://sagemaker-groundtruth-demo/SSDB00005.JPG"}  
{"source-ref": "s3://sagemaker-groundtruth-demo/SSDB00006.JPG"}  
{"source-ref": "s3://sagemaker-groundtruth-demo/SSDB00007.JPG"}
```

“source-ref” で S3 のパスを指定。

あるいは “source” に直接テキストを書くこともできる

## 2. タスクの定義

- 画像分類
- 物体検出
- テキスト分類
- セマンティック  
セグメンテーション

あるいは、

- ユーザ定義のカスタムタスク


**Task type** [Info](#)

**Task selection**  
Select the task that a human worker will perform to label objects in your dataset.

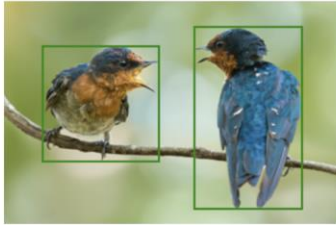
**Image classification**  
Get workers to categorize images into specific classes. [Info](#)

**Basketball**

**Soccer**



**Bounding box**  
Get workers to draw bounding boxes around specified objects in your images. [Info](#)




**Text classification**  
Get workers to categorize text into specific classes. [Info](#)

**Positive**

**Negative**

*'The movie tells a lovely and wise story with honesty and has been acted out with unassuming grace.'*

**Semantic segmentation**  
Get workers to draw pixel level labels around specific objects and segments in your images. [Info](#)





# 3. ワーカーの選択

- Public
  - Amazon Mechanical Turkを利用
- Private
  - 社員などワーカーを自ら調達
- ベンダー
  - 3rd パーティベンダーに依頼
- 追加オプションも選べる
  - 自動ラベリング
  - 複数ワーカーによるラベル付け

## Select workers and configure tool

**Workers** [Info](#)

Worker types

- Public**  
An on-demand 24/7 workforce of over 500,000 independent contractors worldwide powered by Amazon Mechanical Turk.
- Private**  
A team of workers that you have sourced yourself, including you own employees or contractors for handling data that needs to stay within your organization.
- Vendor managed**  
A curated list of third party vendors that specialize in providing data labeling services, available via the AWS Marketplace.

Private teams

Choose from the teams you created in the private workforce or if you need to create a new team, save your progress and go to Labeling workforces to create a new one.

test ▼

► **Additional configuration - optional**  
Automated data labeling, workers per dataset object

## 4. ラベリングツールの設定

- アノテーションの指示書を書く
- 良い例・悪い例を記載
- ラベルを設定

### Image classification labeling tool

Preview 

Use this tool to construct an interface for your workers. Provide labeling instructions with examples that your workers view while performing your labeling task. You can see up to 12 samples of your dataset. Choose a sample and then choose Preview to see your worker's view of the interface.

H1 H2 B I A  

#### Good example

Enter description to explain the correct label to the workers

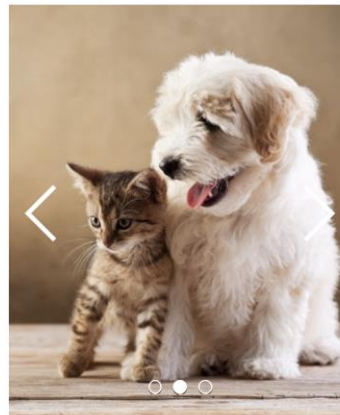


#### Bad example

Enter description of an incorrect label



Image Classification

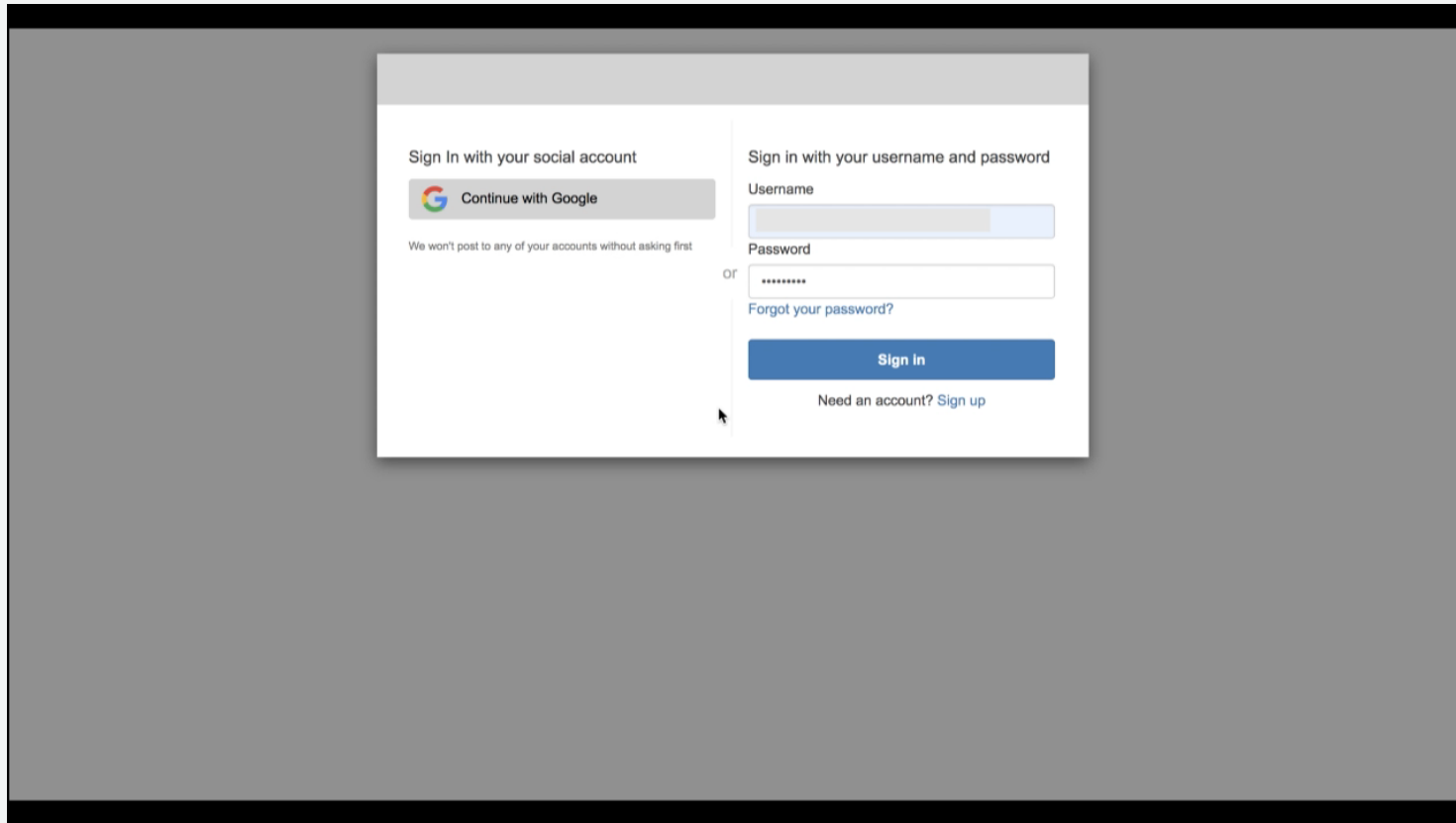


#### Select an option

Add up to 10 labels

Add label

# ワーカー側の利用の流れ



The screenshot displays a sign-in interface with two main sections:

- Sign In with your social account:** Features a "Continue with Google" button. Below it, a note states: "We won't post to any of your accounts without asking first".
- Sign in with your username and password:** Includes a "Username" field, a "Password" field (masked with dots), a "Forgot your password?" link, a blue "Sign in" button, and a "Need an account? Sign up" link.

An "OR" separator is positioned between the two sections. A mouse cursor is visible near the bottom of the form area.

# 出力 (拡張マニフェストファイル)

## output.manifest

```
{ "source-ref": "s3://sagemaker-groundtruth-demo/SSDB00001.JPG",
  "GroundTruthDemo": {
    "annotations": [
      {"class_id": 0, "width": 54, "top": 482, "height": 39, "left": 337},
      {"class_id": 0, "width": 69, "top": 495, "height": 53, "left": 461},
      {"class_id": 0, "width": 52, "top": 482, "height": 41, "left": 523} ],
    "image_size": [{"width": 1280, "depth": 3, "height": 960} ] },
  "GroundTruthDemo-metadata": {
    "job-name": "labeling-job/groundtruthdemo",
    "class-map": {"0": "Car"},
    "human-annotated": "yes",
    "objects": [
      {"confidence": 0.94},
      {"confidence": 0.94},
      {"confidence": 0.94},],
    "creation-date": "2018-11-26T04:01:09.038134",
    "type": "groundtruth/object-detection" } }
```

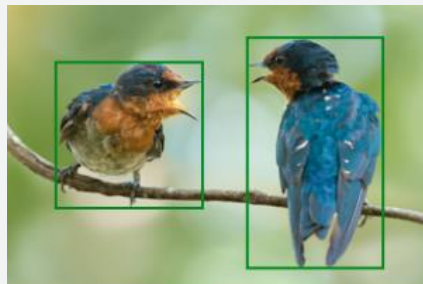
# Demo

# カスタムテンプレート詳細

# 組み込みラベリングツールの利用も 独自実装も可能



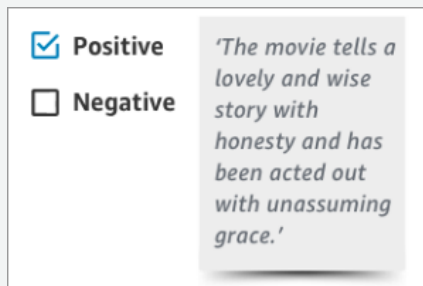
画像分類



物体検出



セマンティック  
セグメンテーション

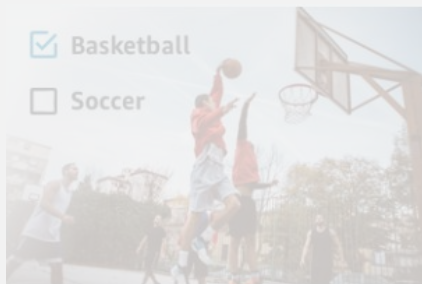


文章分類



カスタム

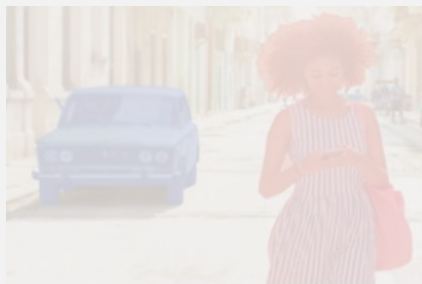
# 組み込みラベリングツールの利用も 独自実装も可能



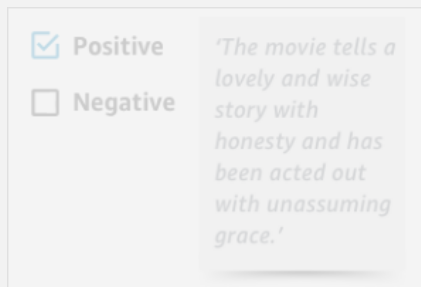
画像分類



物体検出



セマンティック  
セグメンテーション



文章分類



カスタム



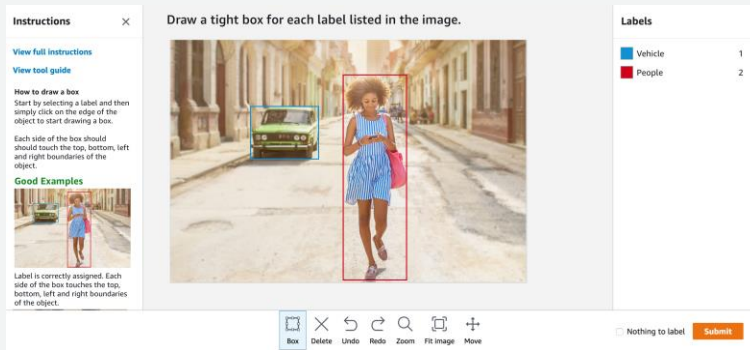
# カスタムテンプレート

前・後処理のLambda関数とラベリングツールのHTMLテンプレートを変更することで、様々なタスクに対応できる機能

前処理用  
Lambda

ラベリング  
HTMLテンプレート

後処理用  
Lambda



Instructions

View full instructions  
View tool guide

How to draw a box  
Start by selecting a label and then simply click on the edge of the object to start drawing a box.

Each side of the box should touch the top, bottom, left and right boundaries of the object.

Good Examples

Label is correctly assigned. Each side of the box touches the top, bottom, left and right boundaries of the object.

Draw a tight box for each label listed in the image.

Labels

Vehicle	1
People	2

Nothing to label Submit

Box Delete Undo Redo Zoom Fit image Move



# カスタムテンプレート

## 前処理用Lambdaの作成

- 入力データの記述されたマニフェストファイルの各項目を読み込み、それをテンプレートエンジンに返す処理を記述

## HTMLテンプレート作成

- テンプレートエンジンのLiquidを採用
- 簡単なサンプルも多数用意されている

## 後処理用Lambda作成

- ワーカーが処理を終了した際の後処理を記述

※ Lambdaはドキュメントのサンプルコードをベースに作成する必要あり  
HTMLテンプレートは多数のサンプルの中から選択できる

# カスタムテンプレート

## HTMLテンプレートのサンプル例

Vision
Bounding Box (prefill)
Bounding Box (multi-class)
Bounding Box (single class)
Facial Detection
Image Classification
Image Contains
Image Moderation
Image Similarity
Image Summarization
Image Tagging
Keypoint
Satellite Image Classification

Semantic Segmentation
Video Classification (from URL)
Video Classification (from Youtube)
Language
Audio Naturalness Evaluation
Audio Transcription
Collect Utterance
Conversation Relevance
Document Classification
Emotion Detection (text)
Emotion Detection (audio)
Fill in the Blank
Fill out the sentence

Same Speaker Evaluation
Semantic Similarity
Sentiment Analysis (text)
Sentiment Analysis (audio)
Sentiment Analysis (tweet)
Translation Quality
Other
<b>Custom</b>
Item Equality
Search Relevance
Survey Link
Website Classification
Website Collection

# カスタムテンプレート

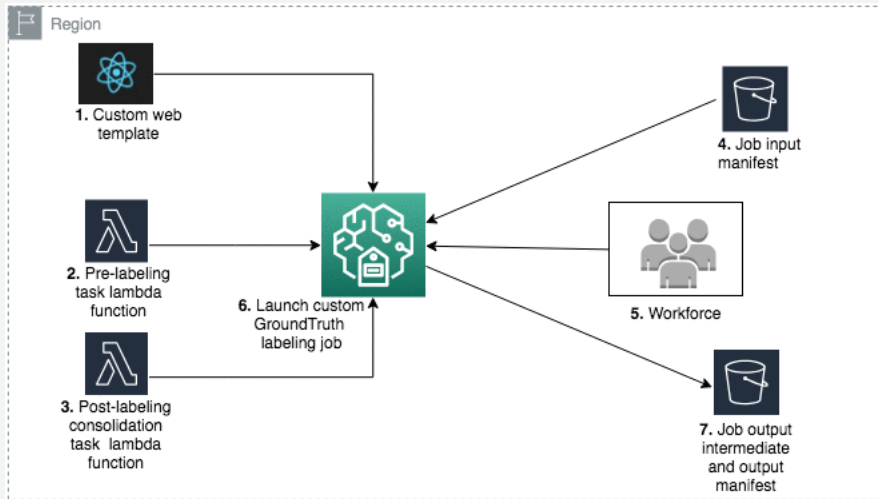
## アノテーション対象のデータ形式

- 画像、映像、音声、文章など様々なファイルをUIに表示できる
- 複数のデータを扱うことも可能（2画像の類似度推定等）

## ラベルの形式

- 既存テンプレートにあるラベル
  - 多クラス分類
  - 画像のセマンティックセグメンテーション
  - 画像の物体検出（カスタムテンプレートであれば複数クラスも可能）
- テキスト（画像に説明文を付与するなど）
- 詳細はドキュメントに

# Amazon SageMaker Ground Truth を使ったカスタムデータラベリングワークフローの構築



## Instructions

The task can be completed with blank, or saved and returned to when time is available to make more progress. If there is evidence in the record to support or deny abstract quality, **highlight** it with the cursor and select **Yes** or **No**. Add any notes you have for each task in the **Notes** free text area.

Min. Size: 81, Actual: Size: 1000, 1 / 17 (200%) Printed: 4 January 2018 AWS 2025-01-04 08:43:21

29 Dec 2017  
aws://1801.0000661/astro-ph.GA/

### Exploring the nature and synchronicity of early cluster formation in the Large Magellanic Cloud: III. Horizontal Branch Morphology

R. Wagner-Kaiser<sup>1</sup>, Dougal Mackey<sup>2</sup>, Ata Sarajedini<sup>1,3</sup>, Roger E. Cohen<sup>4</sup>, Dong Geisler<sup>5</sup>, Soung-Chul Yang<sup>6</sup>, Aaron J. Grocholski<sup>7</sup>, Jeffrey D. Cummings<sup>8</sup>

<sup>1</sup>University of Florida, Department of Astronomy, 215 Physics Drive, Gainesville, FL 32611, USA  
<sup>2</sup>University of Central Florida, Department of Astronomy, 16806 University Blvd, Orlando, FL 32816, USA  
<sup>3</sup>Florida Atlantic University, Department of Physics, 177 Glades Rd., Boca Raton, FL 33432, USA  
<sup>4</sup>Space Telescope Science Institute, Baltimore, MD 21218, USA  
<sup>5</sup>Department of Astronomy, University of Colorado, Campus 300 C, Coonsumo, Ohio Center for Astrophysics, Swanton, Ohio 43884, USA  
<sup>6</sup>Department of Physics and Astronomy, University of Virginia, Charlottesville, VA 22904, USA  
<sup>7</sup>Department of Physics and Astronomy, Swarthmore College, Swarthmore, PA 19081, USA  
<sup>8</sup>Center for Astrophysical Sciences, Johns Hopkins University, Baltimore, MD 21218

CONCLUSIONS

Astronomy Swarthmore College, Swarthmore, PA 19081, USA Center for Astrophysical Sciences, Johns Hopkins University, Baltimore MD 21218  
**ABSTRACT** We leverage new high-quality data from Hubble Space Telescope program GO-14164 to explore the variation in horizontal branch morphology among globular clusters in the Large Magellanic Cloud (LMC). **BACKGROUND:** Our new observations lead to photometry with a precision commensurate with that available for the Galactic globular cluster population. Our analysis indicates that, even metallicity is accounted for, clusters in the LMC largely share similar horizontal branch morphologies regardless of their location within the system. Furthermore, the LMC clusters possess on average, slightly redder morphologies than most of the inner halo Galactic population; we find, instead, that their characteristics tend to be more similar to those exhibited by clusters in the outer Galactic halo. Our results are consistent with previous studies showing a correlation between horizontal branch morphology and age.

**Key words:** (galaxies) Magellanic Clouds, pulsars: star clusters: general, (Galaxy:) globular clusters: general, stars: horizontal branch

**INTRODUCTION**

The variation in horizontal branch (HB) morphology among the Galactic system of globular clusters is known to be strongly, but not strictly, dominated by metallicity. Early studies found a clear distinction between the metal-rich GCs, which generally have very red HBs, and the metal-poor GCs, which tend to be largely populated by the blue side of the red clump (Meyers 1975; Jones & Seargeant 1982; Seargeant 1982). However, this trend is not universal and there are a number of exceptions, particularly in the intermediate-metallicity range of the Milky Way globular clusters. An early example was found by Hodge & Whitman (1985) in studying seven GCs (e.g. M3, M55) that have blue than expected HB morphologies despite their intermediate metallicity. Essentially, metallicity is not alone sufficient to explain HB morphology and additional factors are necessary. This is more succinctly referred to as the “metallicity-metals” effect (Hodge & Whitman 1986, see also Bergh 1996).

Early suggestions to explain the metallicity-metals effect in HB morphology were cluster-to-cluster variations in age and/or helium abundance (see also Bergh 1996, 1997). Other explanations, among which, variations in metallicity in the progenitor stars, in addition to age and helium, were also suggested (see, for example, Sarajedini et al. 2007; Sarajedini et al. 2008; Sarajedini et al. 2009; Sarajedini et al. 2010; Sarajedini et al. 2011; Sarajedini et al. 2012; Sarajedini et al. 2013; Sarajedini et al. 2014; Sarajedini et al. 2015; Sarajedini et al. 2016; Sarajedini et al. 2017; Sarajedini et al. 2018; Sarajedini et al. 2019; Sarajedini et al. 2020; Sarajedini et al. 2021; Sarajedini et al. 2022; Sarajedini et al. 2023; Sarajedini et al. 2024; Sarajedini et al. 2025).

While the underlying cause of the metallicity-metals effect remains a topic of research, this phenomenon has played a significant role in our comprehension of Galactic formation.

**Notes:**

Missing Limitations

Is this a good Abstract?

Yes

No

Submit

# 自動ラベリング詳細

# 自動ラベリング

データの一部をワーカーがラベル付けするだけで、  
残りのラベル付けが自動化され、時間とコストを大幅に削減



※ 5000データ以上の大規模データセットに対して利用可能なオプション機能

# 自動ラベリングの仕組み

## 一部のデータ

アノテーション  
前データ



ワーカーによる  
アノテーション





# 自動ラベリングの仕組み

アノテーション  
前データ



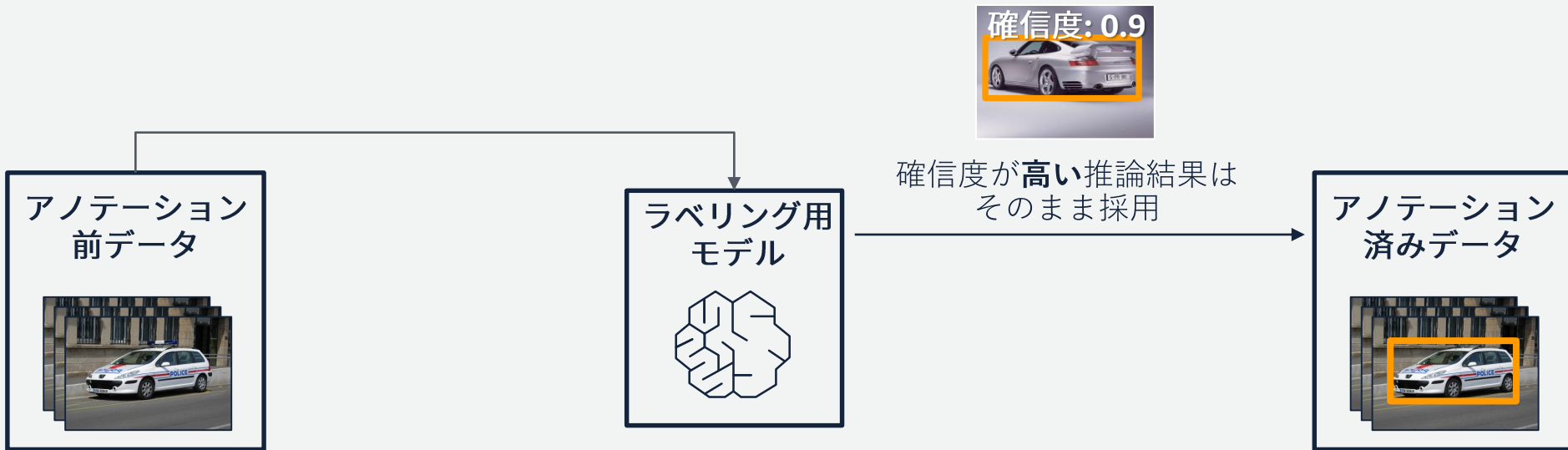
ワーカーによる  
アノテーション



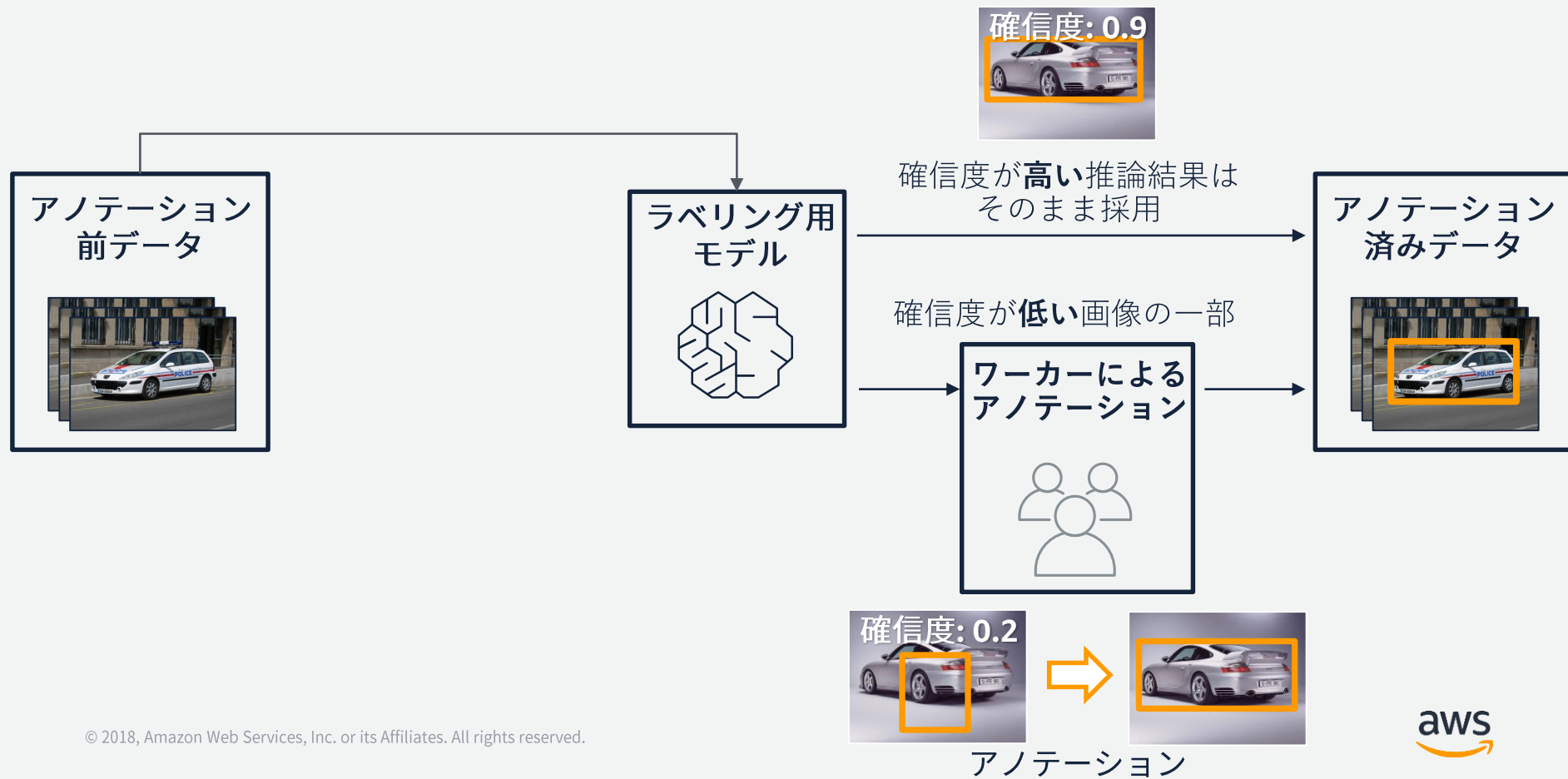
ラベリング用  
モデル



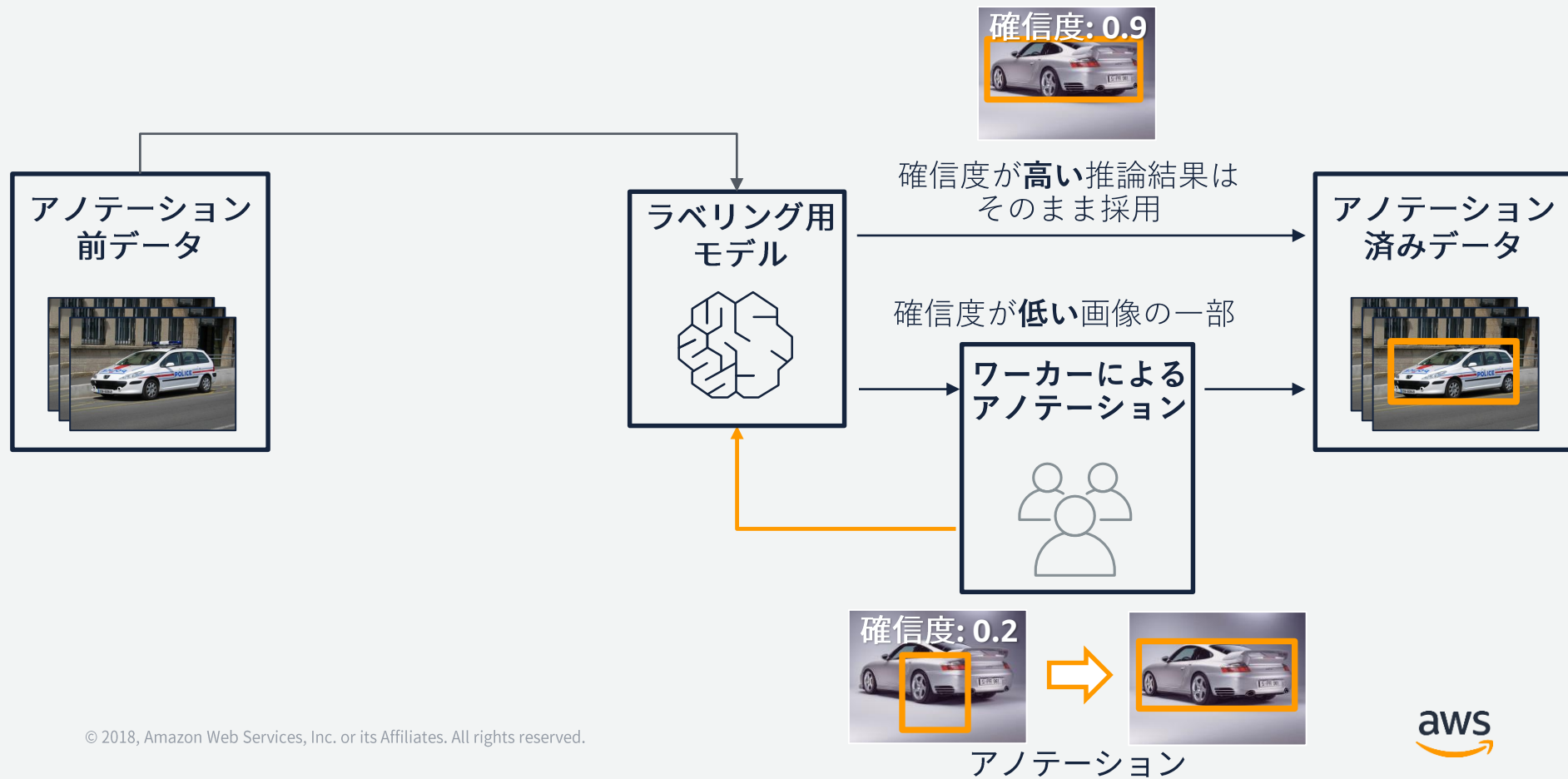
# 自動ラベリングの仕組み



# 自動ラベリングの仕組み



# 自動ラベリングの仕組み



# 価格と提供リージョン

## 価格

- ラベル付けした対象の数に応じた利用料 \$0.08 / 個 (5万個以上はより安価)
- Amazon Mechanical Turk および 外部ベンダを利用する際の利用料
- 自動ラベリング利用時は、裏で動く SageMaker の学習/推論の利用料

## 提供リージョン

- バージニア北部 / オレゴン / オハイオ / アイルランド / **東京**

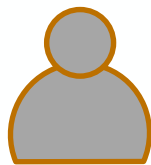


# Appendix

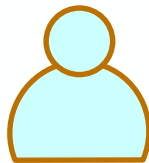
# ラベルの決定方法



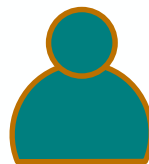
# ラベルの決定 (Label Consolidation): 多数決の場合



bulldog



sharpei



bulldog



bulldog



**bulldog**

# ラベルの決定 (Label Consolidation): 正答率による評価



(犬に詳しい人)



bulldog



sharpei



bulldog



bulldog

正しいラベルを  
選ぶ確率



bulldog 0.1  
sharpei 0.9

- 単純な多数決ではなく、ワーカーの正答率をもとにラベルを決定することでアノテーションの質を担保
- 何人のワーカーの結果をまとめるかはジョブ作成時に指定

# パブリックワーカー詳細

# パブリックなワーカーを利用する

プライベートとの違いは以下の3つのみ

- ジョブ作成時にワーカータイプ > パブリック を選択
- タスクあたりの料金を設定
- 確認項目にチェック

## ワーカーの選択とツールの設定

ワーカー [Info](#)

ワーカータイプ

**パブリック**  
Amazon Mechanical Turk による世界中の 50 万人以上の独立系請負業者のオンデマンド 24 時間年中無休のワーカー。

**プライベート**  
組織内に留まる必要があるデータを処理する従業員や請負業者を含め、自身が調達したワーカーのチーム。

**ベンダー管理下**  
データラベリングサービスの提供を専門とするサードパーティベンダーのリストで、AWS Marketplace から入手できます。

タスクあたりの料金  
支払う価格は、ワーカーがタスクを完了するのにかかる時間によって異なります。

\$0.012  
完了するのに 5 秒以上かかる ▼

このデータセットにはアダルトコンテンツは含まれていません。 [Info](#)

私は自分のデータセットが Amazon Mechanical Turk の一般の従業員によって閲覧されることを理解しています。私は自分のデータセットに個人識別情報 (PII) が含まれていないことを認識しています。 [Info](#)

▶ **追加設定 - オプション**  
自動データラベリング、データセットオブジェクトあたりのワーカー

# パブリックなワーカーを利用する

あとはワーカーがアサインされ、タスクが終了するのを待つだけ

Amazon SageMaker > ラベリングワークフォース

パブリック | プライベート | ベンダー

Amazon Mechanical Turk によって提供される、世界各地のオンデマンドのワーカーからなるチームです。

## パブリックチームの概要

名前 Public_workforce_team	ARN arn:aws:sagemaker:ap-northeast-1:394669845002:workteam/public-crowd/default
-----------------------------	--

## ラベリングジョブ (2)

🔍 名前でラベリングジョブを検索

< 1 > ⚙️

名前 ▼	ステータス ▼	タスクのタイプ ▼	ラベル付きオブジェクト/合計 ▼	作成時刻 ▼
dogcat-image-classification-0414-public	🔄 進行中	イメージの分類	7 / 8	2019年4月14日 8:34 UTC
dogcat-image-classification-0414-clone	🔄 進行中	イメージの分類	8 / 8	2019年4月14日 8:28 UTC

# パブリックなワーカーを利用する

- タスクに対して設定する料金が安すぎると、人が集まらない場合やアノテーションの質が低下する場合あり
- アノテーションの質を向上させるためには、タスクの説明文章の改善や、1データにアノテーションするワーカーの数を増やすことで対応

# その他

# Amazon SageMaker Ground Truth と自動化されたデータのラベル付けによる低コストでのデータのアンノテーション

Amazon Web Services ブログ

## Amazon SageMaker Ground Truth と自動化されたデータのラベル付けによる低コストでのデータのアンノテーション

by Krzysztof Chalupka | on 07 FEB 2019 | in SageMaker | Permalink | Share

Amazon SageMaker Ground Truth を使うと、正確にラベル付けされた機械学習データセットを簡単に低価格で構築することができます。ラベル付けのコストを削減するために、Ground Truth の機械学習を使用して、人によるアンノテーションが必要な「困難な」画像と、機械学習で自動的にラベル付けできる「簡単な」画像を選択します。この記事では、自動化されたデータのラベル付けの仕組みと、その結果の評価方法について説明します。

### 自動化されたデータのラベル付けを伴う物体検出ジョブを実行する

以前のブログ記事では、Julien Simon が AWS マネジメントコンソールを使ってデータのラベル付けジョブを実行する方法を説明しました。このプロセスをより軽く制御するには、API を使用できます。その方法をご紹介しますため、今回は鳥の画像 1,000 個に対して **バウンディングボックスアンノテーション** を生成する API を使用した Amazon SageMaker Jupyter ノートブックを使用します。

注意: デモノートブックの実行コストは約 200 USD です。

デモノートブックにアクセスするには、ml.m4.xlarge インスタンスタイプを使用して Amazon SageMaker ノートブックインスタンスを開始します。インスタンスは、このステップバイステップチュートリアルに従ってセットアップできます。ステップ 3 では、IAM ロールの作成時に「任意の S3 バケット」にチェックを入れるようにしてください。以下にあるように、Jupyter ノートブックを開いて [SageMaker Examples] タブを選択し、object\_detection\_tutorial.ipynb を起動します。

The image shows two side-by-side screenshots. The left screenshot is the AWS SageMaker console. It displays the 'Notebook instances' page with a table of instances. One instance named 'demo' is in 'InService' status. A red circle highlights the 'Open Jupyter' button in the 'Actions' column. A red arrow points to this button with the text 'In Amazon SageMaker, open any notebook instance.' The right screenshot shows the Jupyter notebook interface. The 'SageMaker Examples' tab is selected and circled in red. Below the tab, there is a list of example notebooks, with the instruction 'Choose the SageMaker Examples tab.' highlighted in red.

<https://aws.amazon.com/jp/blogs/news/annotate-data-for-less-with-amazon-sagemaker-ground-truth-and-automated-data-labeling/>



# Amazon SageMaker Ground Truth を使用して階層型ラベル分類法を作成する

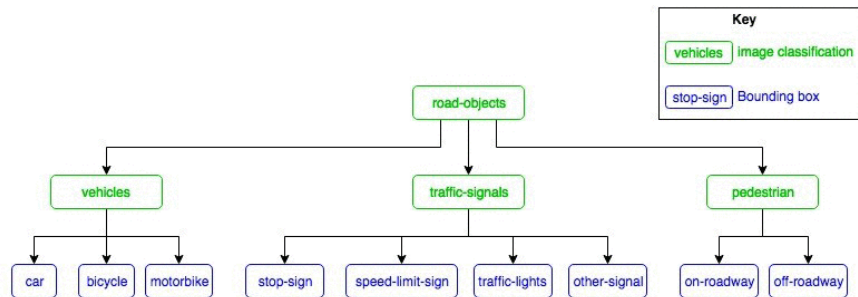
Amazon Web Services ブログ

## Amazon SageMaker Ground Truth を使用して階層型ラベル分類法を作成する

by Doug Youd | on 19 FEB 2019 | in SageMaker | Permalink | Share

re:Invent 2018 で、私たちは [Amazon SageMaker Ground Truth](#) を発表しました。これは、機械学習を使用して [非常に正確なデータセットを構築し、ラベル付けのコストを最大 70% 削減](#) することができます。Amazon SageMaker Ground Truth を使用すると、パブリックおよびプライベートでラベル付けを行う人間の作業者に簡単なアクセスと、一般的なラベル付けタスクのための組み込みのワークフローとインターフェースが提供されます。さらに、Amazon SageMaker Ground Truth は自動データラベル付けを使用してラベル付けのコストを削減します。自動データラベル付けは、人間がラベルを付けたデータから Ground Truth をトレーニングし、サービスが独自にデータにラベルを付けることを学習することによって機能します。

路上でカメラで撮影した画像の大規模なコーパスがあるとしましょう。それぞれの画像には、無人自動車用のアルゴリズムを開発するために重要な多くの異なる対象物 (たとえば、車または交通標識) が含まれている可能性があります。最初に、画像から取得したい情報の階層表現を定義する必要があります (そのようなラベル分類法がどのように見えるかの例については、下記を参照してください)。次に、これらのラベルが付いていない未処理の画像を取得し、高レベルのクラス (「車」、「交通標識」、「歩行者」など) でラベル付けすることによってラベル付けプロセスを開始します。



このブログ記事では、ジョブをチェーンさせて拡張マニフェスト機能を利用することで、Amazon SageMaker Ground Truth を使用してこのような階層的なラベル付けを実現する方法を紹介します。

<https://aws.amazon.com/jp/blogs/news/creating-hierarchical-label-taxonomies-using-amazon-sagemaker-ground-truth/>

# Amazon SageMaker Ground Truth でのラベリングジョブ用の優れた説明の作成

Amazon Web Services ブログ

## Amazon SageMaker Ground Truth でのラベリングジョブ用の優れた説明の作成

by Tristan McKinney | on 06 APR 2019 | in SageMaker | Permalink | Share

[Amazon SageMaker Ground Truth](#)は、機械学習 (ML) 用の高精度なトレーニングデータセットをすばやく構築するお手伝いをします。ご自身のワークフォース、データラベリングに特化したベンダー管理ワークフォースの選択、または Amazon Mechanical Turk が提供するパブリックワークフォースを使用して、人が生成するラベルを提供することができます。質の高いラベルを取得するには、特にパブリックワークフォースを使用している場合、簡単かつ簡潔で明確な説明が必要です。良い説明が書ければ、アノテーションの品質を向上させることができます。正しく行えば、この作業に時間を費やす価値があります。

このブログ記事では、パブリックワークフォースに効果的な説明を作成するためのベストプラクティスをご紹介します。ここで重要なポイントが 2 点あります。ワークフォースへの認知負荷をできるだけ減らすこと、そして説明を微調整して後で発生する問題を避けるためにもプロセスの早い段階で実験することです。たとえば実験で、データの一部に自分でラベルを付けたら、プロセス全体の中でも小規模なジョブをパブリックワークフォースに行ってもらうことができます。

以下のスクリーンショットは、ワーカーの観点から見て適切な説明のある Ground Truth のバウンディングボックスのラベリングタスク例を示しています。このタスク例では、[Google Open Images Dataset](#) から取得したイメージにある花の周りを囲む四角形の枠を描くようワーカーに伝えます。ワーカーがアノテーションを付けている作業中、イメージの左側にあるサイドバーには短い説明が表示されます。はっきりと要領を得た、かつタスクに特化した説明で、サンプルのイメージに焦点を当てています。

**Instructions** × Draw a box around each flower.

View full instructions  
View tool guide

**Good example**

Each box should be as small as possible.

**Bad example**

Each flower should have one box.

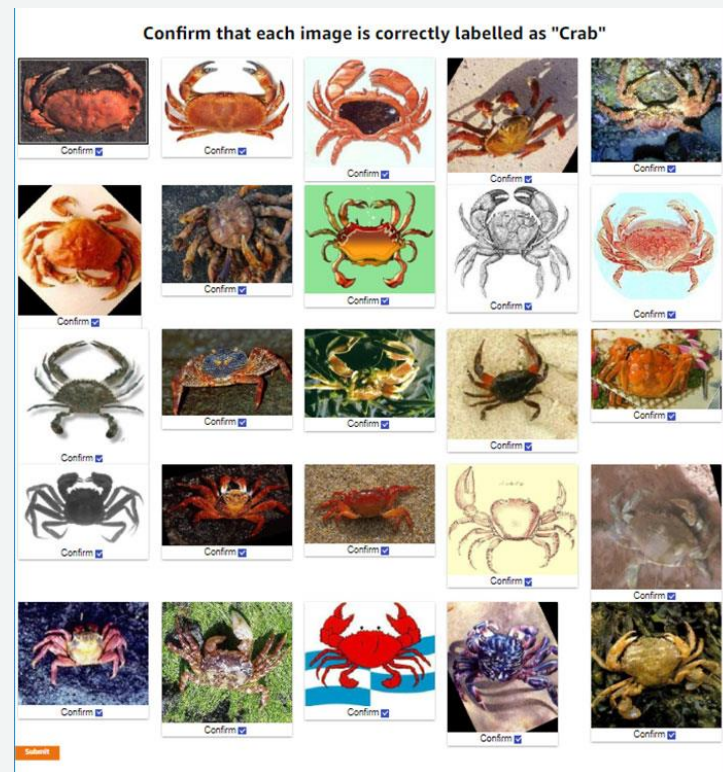
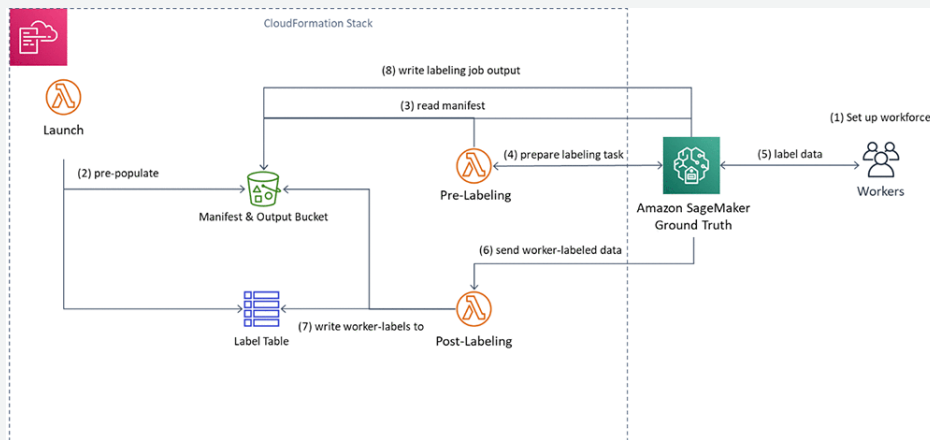
**Labels**

Flower 1

Nothing to label Submit

<https://aws.amazon.com/jp/blogs/news/create-high-quality-instructions-for-amazon-sagemaker-ground-truth-labeling-jobs/>

# Amazon SageMaker Ground Truth を使って大量ラベル付けの品質保証を簡単に行う



# Amazon SageMaker Ground Truth – A Deep Dive with an Interactive Workshop to Build High-Quality Training Datasets

## Take an example: How can we improve

- Task "Draw a bounding box" makes people think "Draw bounding boxes around objects of the specified class in this image."
- It is better if the instructions have an image similar to the annotated ones.
- Panel instruction is too long and should be shorted



## Object detection



- Common solution: Use label from a single, trusted worker
- Our approach: Find corresponding boxes based on overlap, and average the coordinates, favoring smaller boxes. Reject unmatched boxes