

DEV DAY

20-22.10.2020

F - 6

Tensorflow on Lambda ~機械学習予測処理の低コスト化~

中川 帝人
上席課長/シニアデータサイエンティスト
株式会社オープンハウス/情報システム部 ディスラプティブ技術推進G

自己紹介

- 氏名：中川帝人(@TeitoNakagawa)
- 所属：株式会社オープンハウス
- 役職：シニアデータサイエンティスト/上席課長/宅地建物取引士
- 得意分野：データベース、GIS、データ分析・機械学習、AIの企画・開発
- 経歴：日本オラクル、ソフトバンク等で分析コンサルタントから分析・新規事業企画系のマネジメントを担当。2018年3月よりオープンハウスでAI・データ活用の研究開発と活用のチームを牽引。不動産×AI・MLのソリューションを複数実現。
- 好きなAWSのサービス：Lambda



まとめ

まとめ

不規則な社内予測サービス はLambdaで実行する

Usecase

- : 社内システム、不規則なアクセス
- X : コンシューマ向け、規則的なアクセス

Good

- ・インフラ考慮不要
- ・スケールを待たない
- ・レスポンスタイムも早い
- ・同時実行数も多い

Risk

- ・環境構築の検証
- ・モデルサイズ上限（EFSを利用）



Agenda

1. オープンハウスにおけるAI活用
2. オビ付け業務の自動化
3. 課題：不規則かつ急激に上がる負荷
4. Tensorflow on Lamda
5. 検証結果



オープンハウスにおけるAI活用

会社紹介

<https://oh.openhouse-group.com/kanto/about/>

•株式会社オープンハウス

- 総合不動産会社
- 戸建開発と販売が主力事業
- 営業中心の会社

•情報システム部

- 内製主義＝コード力重視
- ほぼ全てのシステムでOn Cloud
- 最近Pythonが主流に
- AI・RPAで事例多数あり

好立地、ぞくぞく。

「都心に戸建てを持ちたい」という夢を叶えるために、時代が変わっても、価値が変わらないものをコンセプトに、都心の「好立地の家」を、お届けしています。

オープンハウスは創業22周年を迎える、総合不動産グループです。

オープンハウスは東京23区を中心に横浜・川崎・埼玉・千葉・名古屋・福岡で土地や戸建てを仲介する不動産会社です。そして、企画・開発、建設を担うオープンハウス・ディベロップメントと建築を請け負うオープンハウス・アーキテクトが力を合わせて、“等身大”の家づくりを実現しています。

33.1% 過去最高の売上高・利益
7期連続で更新し続けています。

OPEN HOUSE GROUP 年平均成長率 33.1%
※2013～2020年

年	売上高 (億円)	経常利益 (億円)
2013年	969	91
2014年	1,121	128
2015年	1,793	202
2016年	2,472	291
2017年	3,046	361
2018年	3,907	460
2019年	5,403	549
2020年	5,700	770

オープンハウスは2013年9月の東証一部上場以来、年平均成長率33.1%という高い水準で成長を続けています。創業22周年の2019年には5,000億円を突破。2020年に向けて、さらに高い業績予想を掲げています。

首都圏から名古屋まで
オープンハウスは東京だけじゃない。

オープンハウス・ディベロップメントは、多数商圏で「地域ビルダーNo.1」を獲得
～「19年度全国No.1ホームビルダー大賞集より」～

東京都No.1 千葉県No.1 名古屋市No.1 その他多数エリア多数

仕入 企画・開発 建築 土地・戸建仲介 建築請負

OPEN HOUSE Development OPEN HOUSE Architect

お客様

オープンハウスは東京23区・横浜・川崎・埼玉・千葉を中心に数十店舗を展開。オープンハウス・ディベロップメントは東京・神奈川の多数の地区で着



不動産×AIはニーズ大

- ニーズないでしょ
⇒あります
- 価格予測でしょ？
⇒ニーズない
- 実は結構できることある
 - 超アナログ
⇒ITの活躍の余地が大きい
 - 書類がとにかく多い
⇒AI・MLの活躍の余地が大きい
 - 単純な作業が多い
⇒現状のAIでも十分実現可能

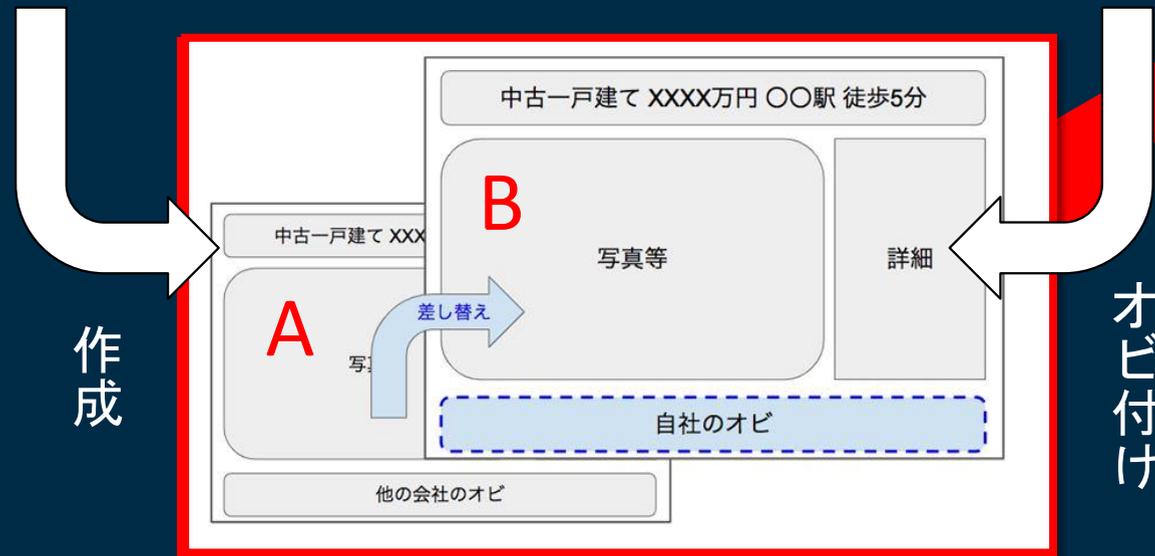


オビ付け業務の自動化

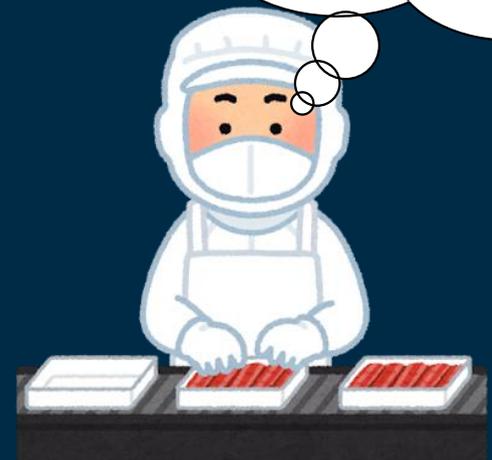
オビ付けとは？

マイソク（＝チラシ）の他社のオビ（＝情報）を
自社のオビに差し替える作業

仲介の流れ



貼り付けるだけの
単純な作業だけど
数が多い



⇒画像処理で実現

物体検知とは

オビ付けは画像処理＝物体検知モデルで実装可能

- 実装はオープンソース
- サンプルはネットに腐るほどある
- 有名なソフトウェア
 - YOLO
 - Tensorflow Object Detection API
- モデル
 - SSD
 - Faster R-CNN

TensorFlow Object Detection API

TensorFlow 2.2 TensorFlow 1.15 Python 3.6

Creating accurate machine learning models capable of localizing and identifying multiple objects in a single image remains a core challenge in computer vision. The TensorFlow Object Detection API is an open source framework built on top of TensorFlow that makes it easy to construct, train and deploy object detection models. At Google we've certainly found this codebase to be useful for our computer vision needs, and we hope that you will as well.

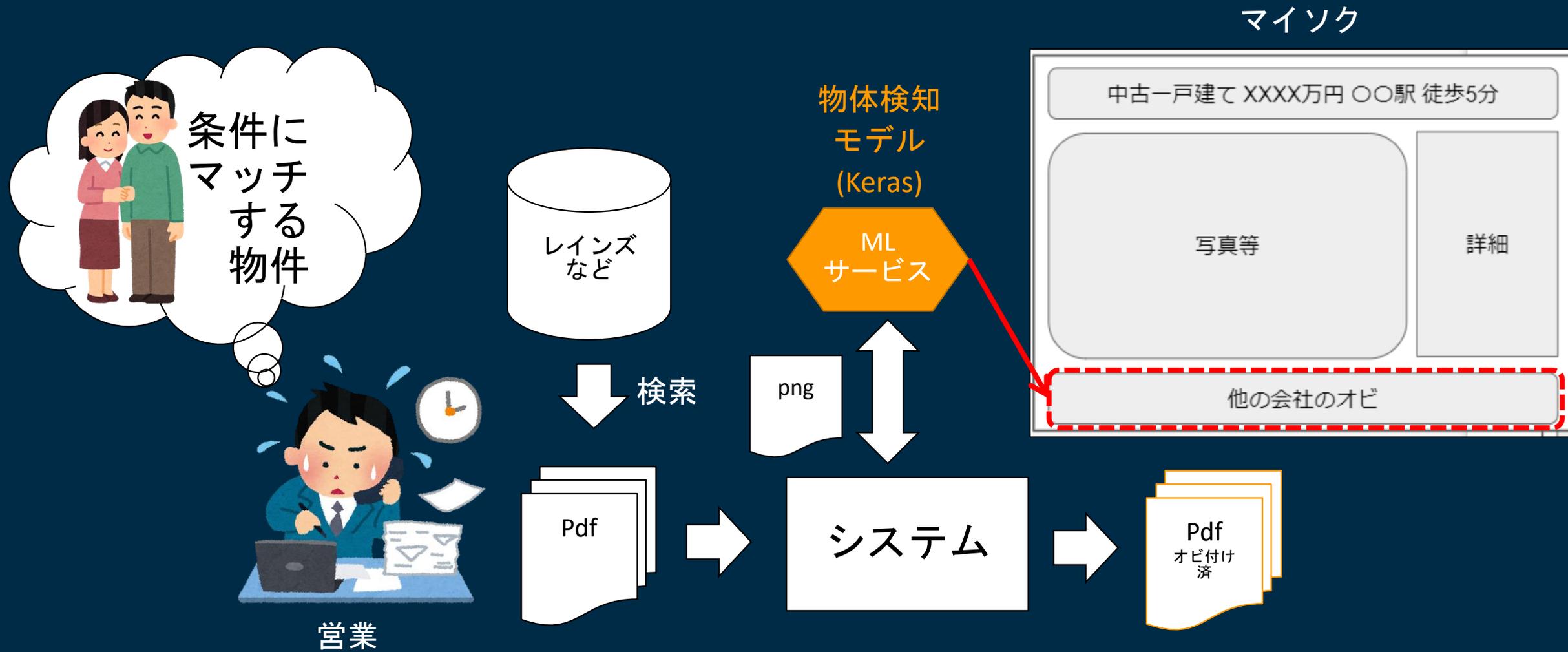


Contributions to the codebase are welcome and we would love to hear back from you if you find this API useful. Finally if you use the TensorFlow Object Detection API for a research publication, please consider citing:

https://github.com/tensorflow/models/tree/master/research/object_detection

作ったシステムの利用イメージ

PDFで送られたマイソクに物体検知モデルでオビを貼り付ける



課題：不規則かつ急激に上昇する負荷

不規則かつ急激に上昇する負荷

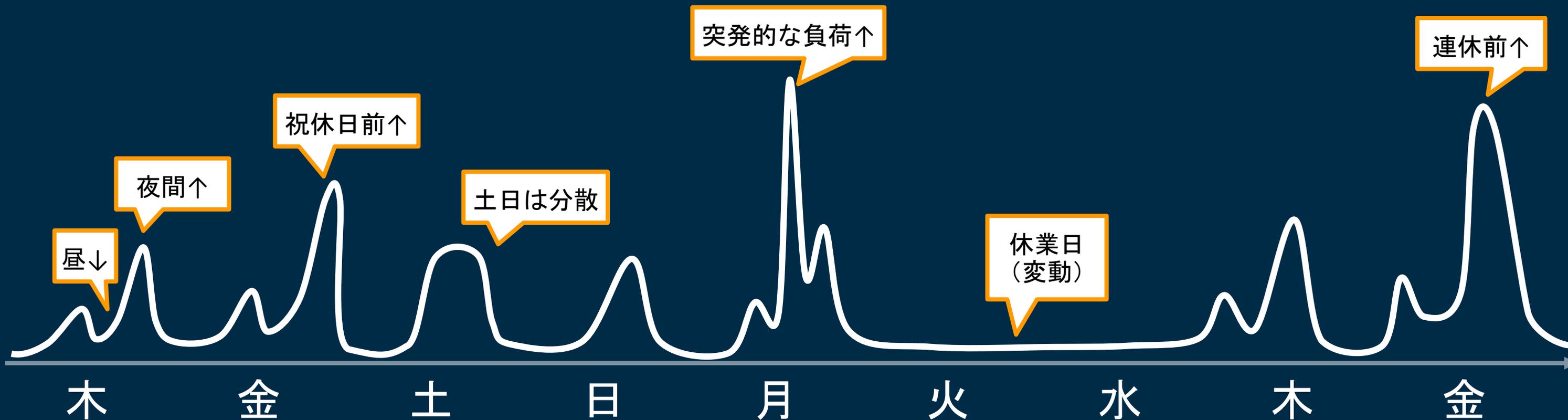
厳しいインフラ要件

1. 1回の予測の負荷が重い：入力画像のサイズを落とすと精度が下がる
2. 1回のリクエストの負荷が重い：一つのpdfが多くて50ページ
3. リクエストが集中：500人近い営業が大体同じ時間に作業
4. 不規則なリクエスト：休みの前日に負荷が上昇、休業日は不規則
5. 高い性能要件：働き方改革で現場の時間は足りない



実際の負荷イメージ

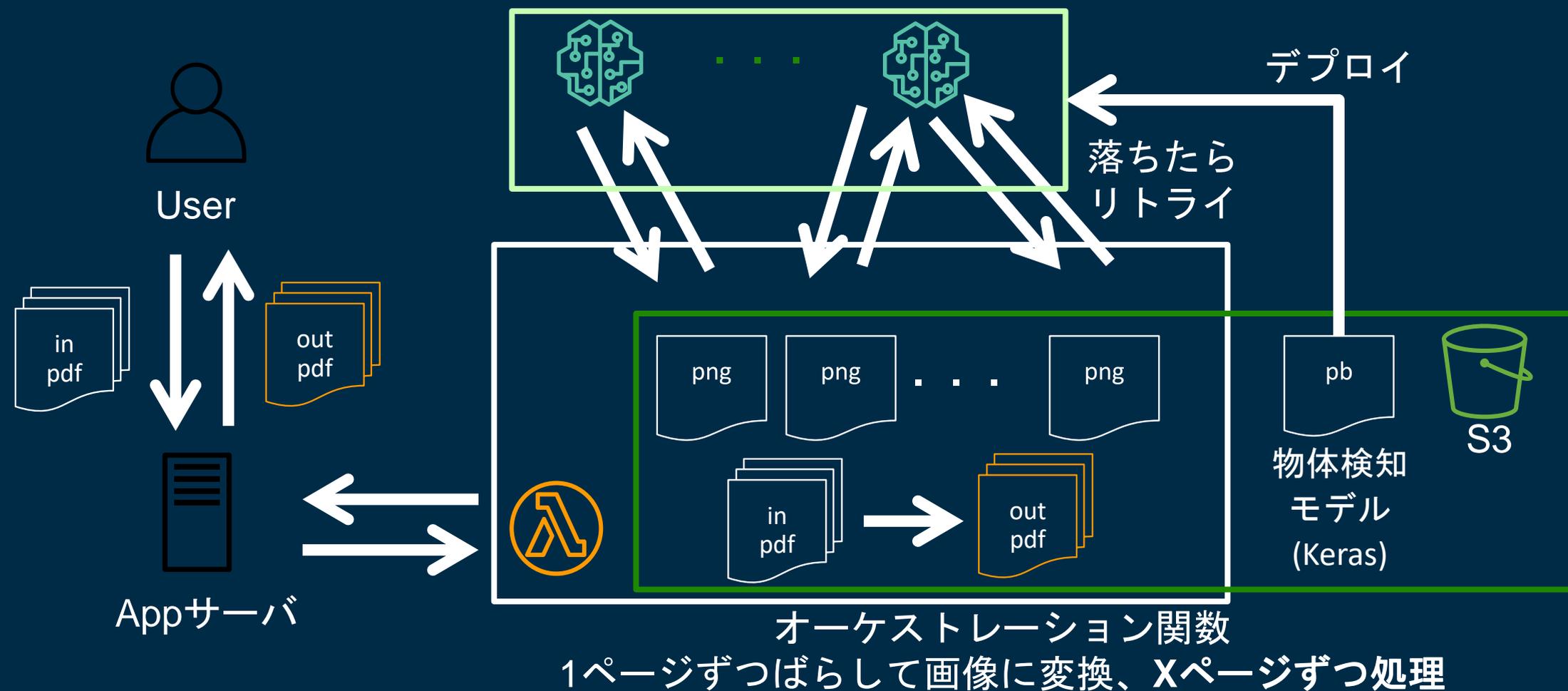
変動要因が多すぎて負荷を予測できない
いつ負荷が上がってもいい形で実装せざるを得ない



アーキテクチャ (ASIS)

MLサービスを複数デプロイ、急激な負荷上昇に耐える

MLサービス：常時X並列のEPで処理を受け付ける。
高負荷の時はそれぞれがスケールする



実装上の課題

- 常時複数サービスを起動した結果、コストが上がってしまった
- 複数サービスでもさばけない事が多々発生し、リトライを実装しなければならなかった
- スケーリングに時間がかかるため無駄なリソースを確保した

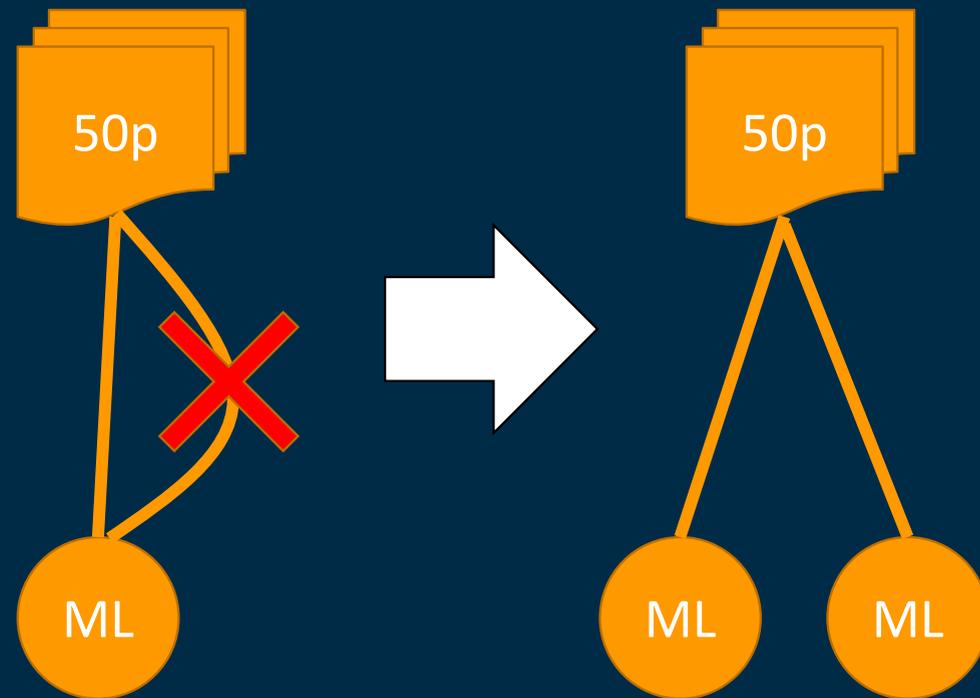
⇒関数で解決できないか？



Tensorflow on Lamda

関数コンピューティングで解決出来るのでは？

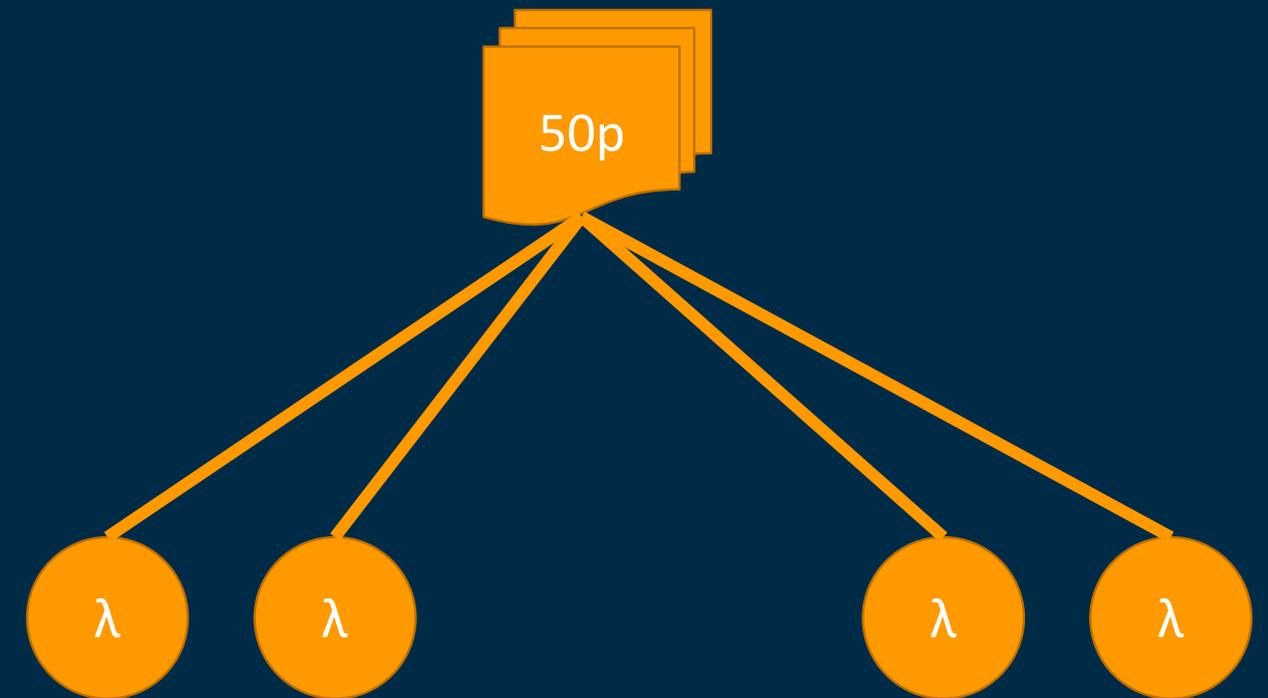
MLサービス



一つのモデルに処理が集中するので落ちる

スケーリングするので落ちない

関数



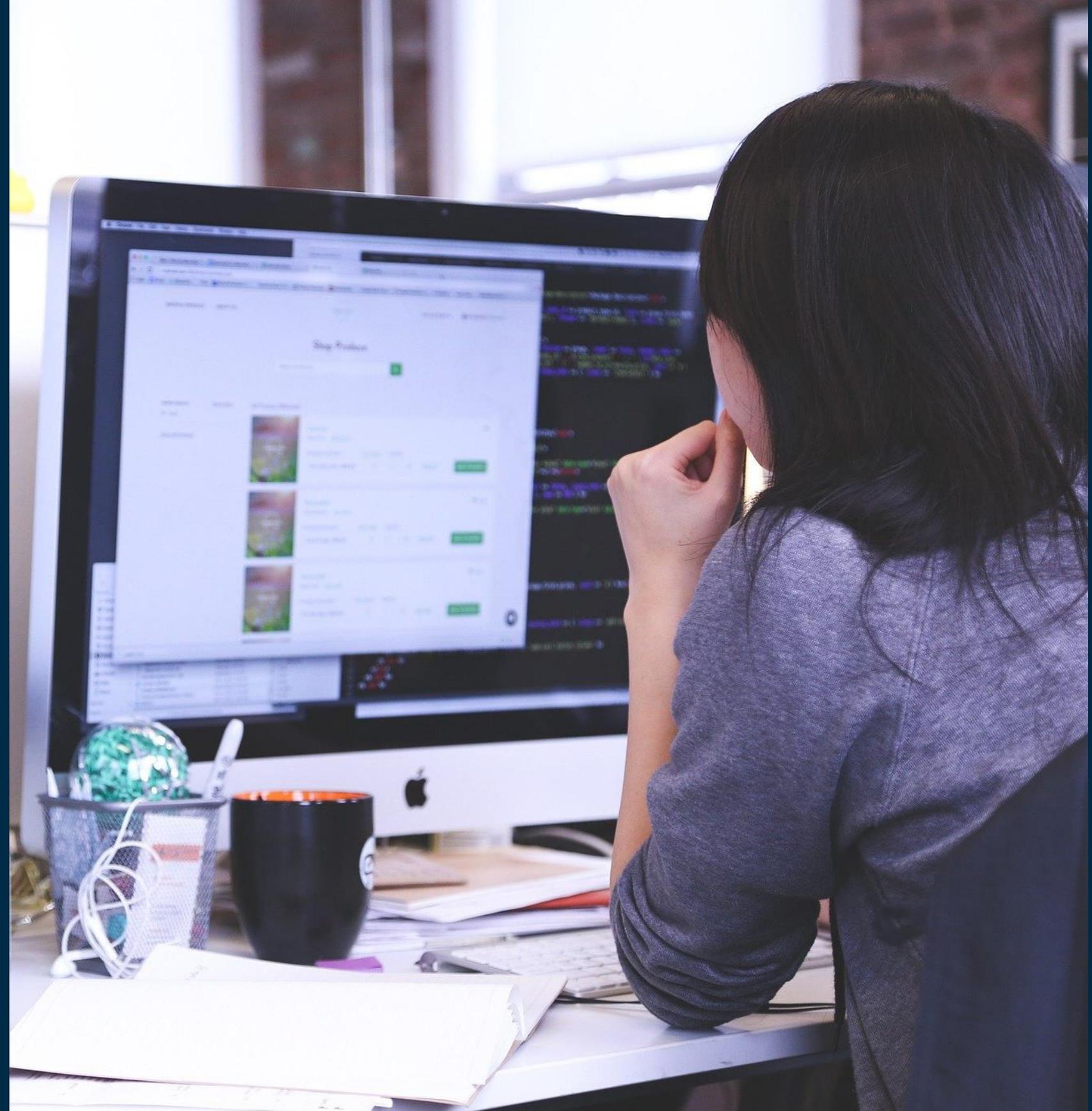
一気に50p並列処理すれば早いのでは？

機械学習予測の方式比較

	サーバ	MLサービス	関数
実装	サーバ(コンテナ)構築 EC2, ECRでスケーリング	Sage Maker	Lambda
開発	独自Appサーバ Tensorflow Serving で関数や環境を構築	Modelを配置する	コードと モジュールのロード
メリット	自由に実装できる	スケーリングは 自動で実行	実行した分だけ課金 不規則なアクセス○
デメリット	スケーリング処理 実装は複雑	スケーリングが遅い	環境の制約が強い 毎回環境が構築される

Lambda実装で工夫したポイント

1. モジュールサイズの削減：
Tensorflow自体が大きいモジュールのため、Serverlessを使ってサイズを削減した
2. 初期化処理をinitに：予測の度に変わらない変数の初期化（環境の構築）は全てinitに記載することで実行時間を節約した



実装概要

- ランタイム : Python 3.6
- メモリ : 3008 MB
- モデル : Keras(Tensorflow)
- モジュール
 - tensorflow=2.0.2,
 - numpy=1.7.2,
 - Pillow=6.1.0,
 - h5py=2.10.0
- Severless : 1.78.1



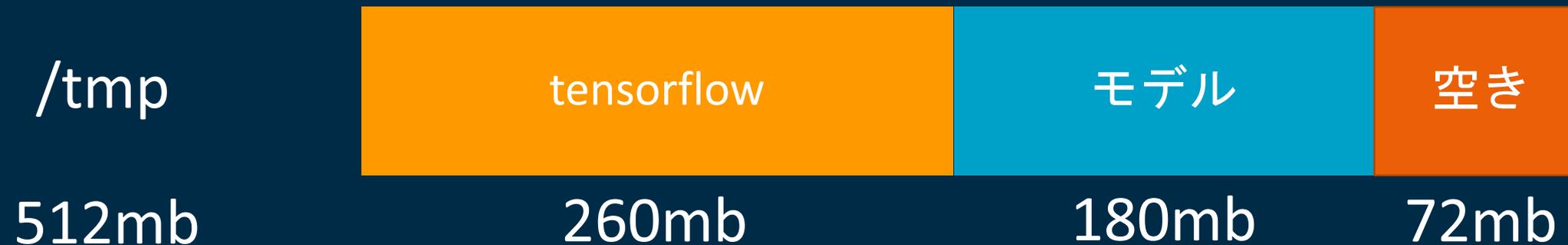
モジュールサイズの削減

Python Lambda最大の問題：パッケージサイズ（特に非ピュアPython）

Serverless Error -----

An error occurred: ClassifyLambdaFunction - Unzipped size must be smaller than 176515274 bytes (Service: AWSLambdaInternal; Status Code: 400; Error Code: InvalidParameterValueException; Request ID: f59a1c34-8226-439d-9d86-04e3f6c5946b; Proxy: null).

ServerlessでTensorflowのディスク利用を削減してデプロイ可能に



Serverless上の設定

Serverless Framework : 関数コンピューティングのためのフレームワーク

- 複数のクラウドサービスに対応
- ローカルでの実行を簡単に

Serverless Python Requirements : Python Requirementsを扱うためのServerlessプラグイン

- dockerizePip : クロスコンパイルでデプロイを高速化
- slim,slimPatterns:不要なモジュールを削除してディスクを減らす

参考 : 「AWS LambdaでTensorflow 2.0を使った画像分類」

https://tech.unifa-e.com/entry/2019/09/17/085400?utm_source=feed



初期化処理をinitに

初期化（モデルダウンロード）に1番時間がかかっていた



処理時間を1/6に削減⇨コストを1/6に削減

コードの実装方法

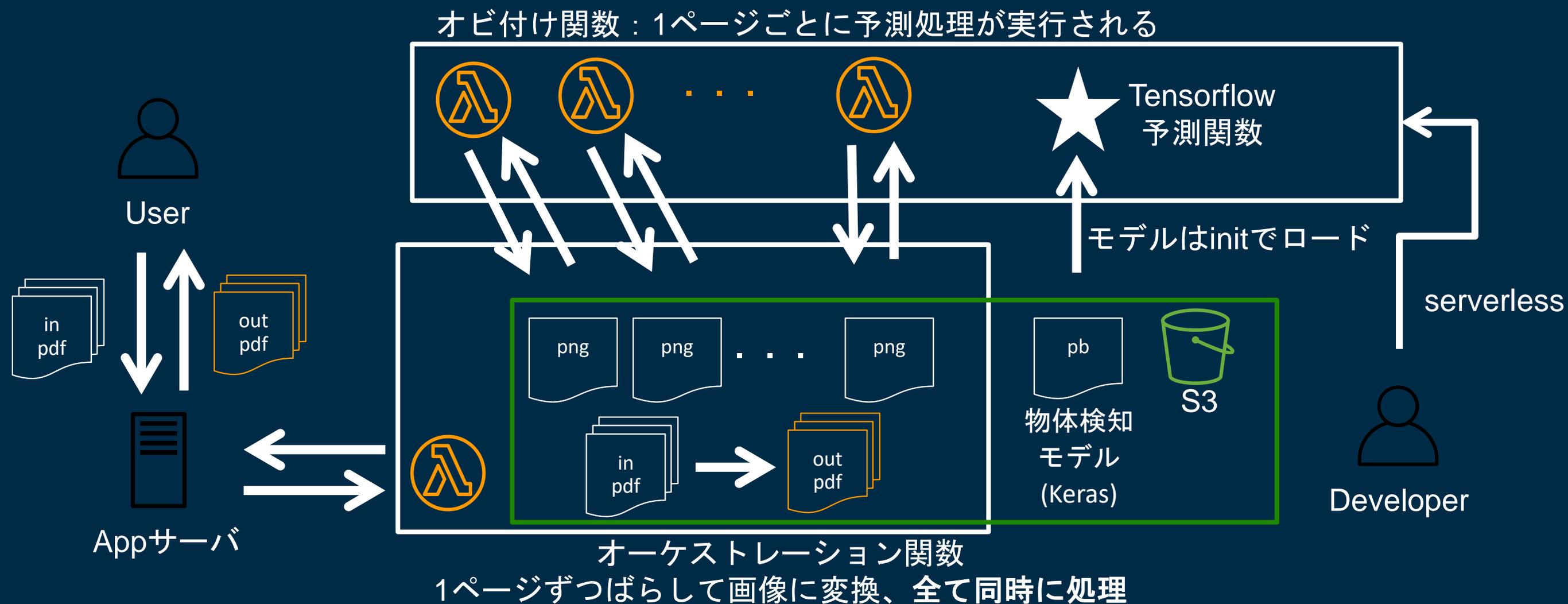
共通した処理はinitに記載すれば1度しか実行されない。

```
...
MODEL_FN = None #予測関数のグローバル変数
...
#初期化：モデルのダウンロード、予測関数の抽出
def init():
    global MODEL_FN
    MODEL_FN = NEW_MODEL.signatures['serving_default']
    ...

#予測：S3画像のダウンロード、前処理、予測
def handler(event, context):
    global MODEL_FN
    ...
    output_dict = MODEL_FN(encoded_image)
    ...
```

アーキテクチャ (TOBE)

ラムダ関数に変更することで全て同時に処理できる



検証結果

検証条件

方法：PythonでローカルPCから
予測関数をコール（非同期）

条件

実行回数：10000回

テスト時間：約33分

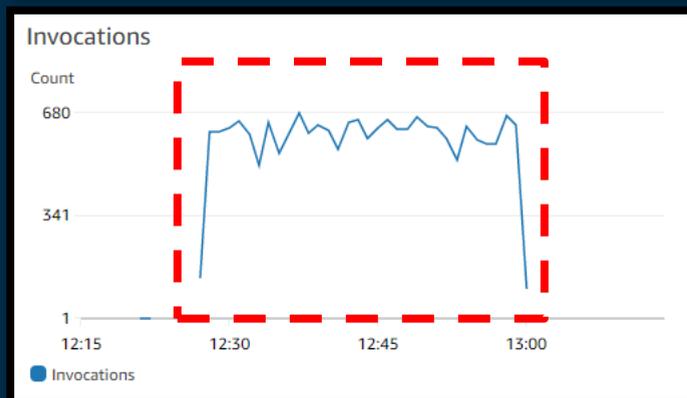
同時実行回数：約285回/分

⇒実際の負荷を大きく超える
条件で検証を実施



検証結果

▶	2020-10-12T12:33:33.361+09:00	START RequestId: 61119c60-40ff-41da-8759-e99d865d414b Version: \$LATEST
▶	2020-10-12T12:33:33.592+09:00	classify
▶	2020-10-12T12:33:36.612+09:00	END RequestId: 61119c60-40ff-41da-8759-e99d865d414b
▶	2020-10-12T12:33:36.612+09:00	REPORT RequestId: 61119c60-40ff-41da-8759-e99d865d414b Duration: 3246.95 ms Billed Duration: 3300 ms Memory Size: 3008 MB Max Memory Used: 2058 MB



•Max Used Memory : **概ね2058MB**

•エラー回数 : 0回

•利用料金 : 約\$1

⇒エラーなしでコストとレスポンスタイムの改善を実現

まとめ

まとめ

不規則な社内予測サービス はLambdaで実行する

Usecase

- : 社内システム、不規則なアクセス
- X : コンシューマ向け、規則的なアクセス

Good

- ・インフラ考慮不要
- ・スケールを待たない
- ・レスポンスタイムも早い
- ・同時実行数も多い

Risk

- ・環境構築の検証
- ・モデルサイズ上限（EFSを利用）

