

SoVeC Smart VideoにおけるAmazon SageMakerを 用いた開発運用事例

Feb. 27th, 2020

ソニー R&Dセンター Tokyo Laboratory 16 大石壮一郎

自己紹介

大石壮一郎

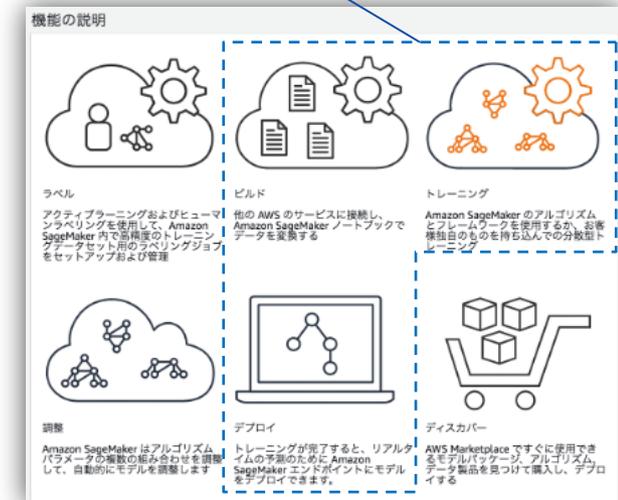
- ✓ 所属: ソニー（株） R&Dセンター
- ✓ 技術: 画像処理、クラウド、機械学習
- ✓ 職種: PLを中心に何でも雑多に
- ✓ 業務: 新規事業の立上げ（通算6回目）
 - ✓ 機械学習関連: 4回、SageMaker: 今回初
 - ✓ 前回: Edgeでの機械学習（Python 2.7）
- ✓ 好きなサービス
AWS Lambda、SageMaker



executive summary (, agenda)

- ✓ SoVeC Smart Videoでは深層学習を用いた超解像技術を提供している
- ✓ 超解像技術はAWS SageMakerの複数機能をベースに組み合わせて構築した
- ✓ SageMakerにより短期間で機械学習のDevOpsを実現した

SoVeC Smart Videoに採用



AWS SageMaker機能群

SoVeC株式会社

SoVeC

ソニーネットワークコミュニケーションズと総合PR会社ベクトルの ジョイントベンチャー

SONY

Sony Network Communications

AIやIoTなど、ソニーグループの
技術を活用したソニーネットワーク
コミュニケーションズの先進的な
サービス事業における知見



vector 

ベクトルグループの
コミュニケーション戦略
実行力を活用

デジタルコミュニケーションの民主化を目指し、日々高度
化する企業のマーケティング活動の最適化・効果の最大化
を図る新たなソリューションを提供

SoVeC Smart Video

コンセプト

SoVeC

デザイナークオリティの動画を、誰でも簡単・リーズナブルに
～デザイナーが伴走するように動画が出来る～



<<https://smartvideo.ovec.net/cp/>>

5

SoVeC Smart Video: サービス概要

SoVeC Smart Video 概要

SoVeC

AIを活用した動画自動生成クラウドサービス

ブラウザだけでプロのデザイナー級の動画を誰でもあっという間に

貴社でお持ちの情報や資料



クラウド上で簡単に動画制作

映像・画像
テキスト

SoVeC Smart Video

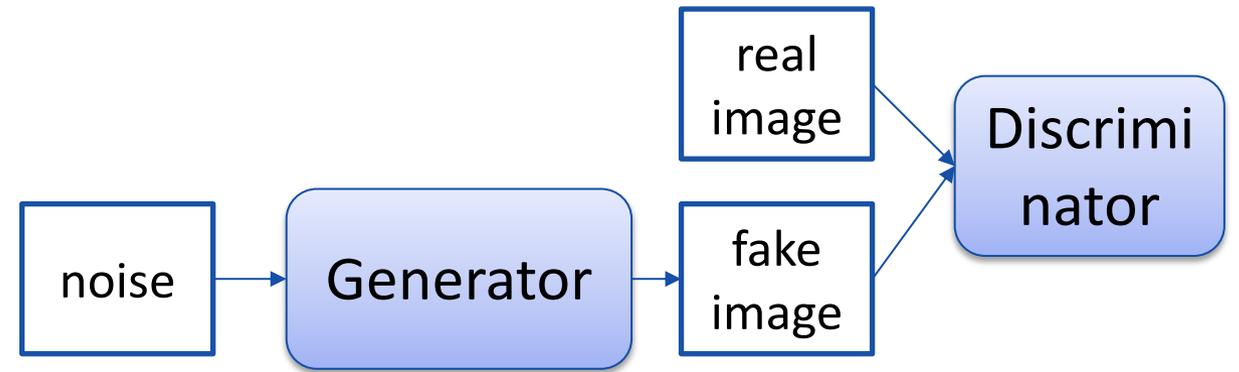
素材を投入



様々な用途に使える動画を出力

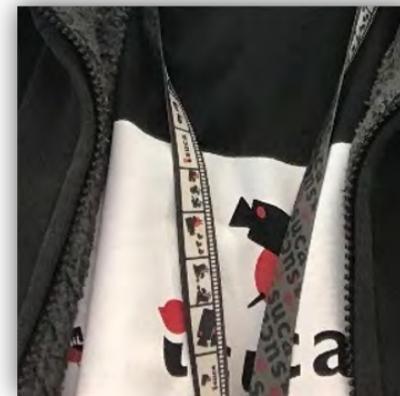


SoVeC Smart Video: SRGAN を用いた超解像機能

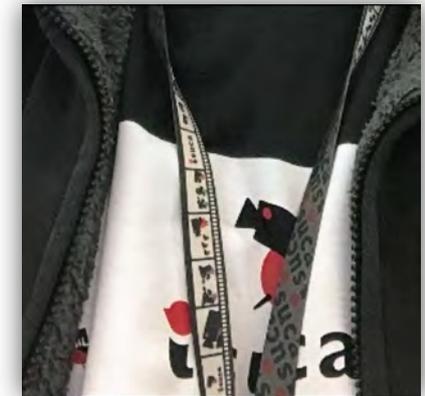


Super Resolution Generative Adversarial Network

低解像な静止画素材を高精細化する

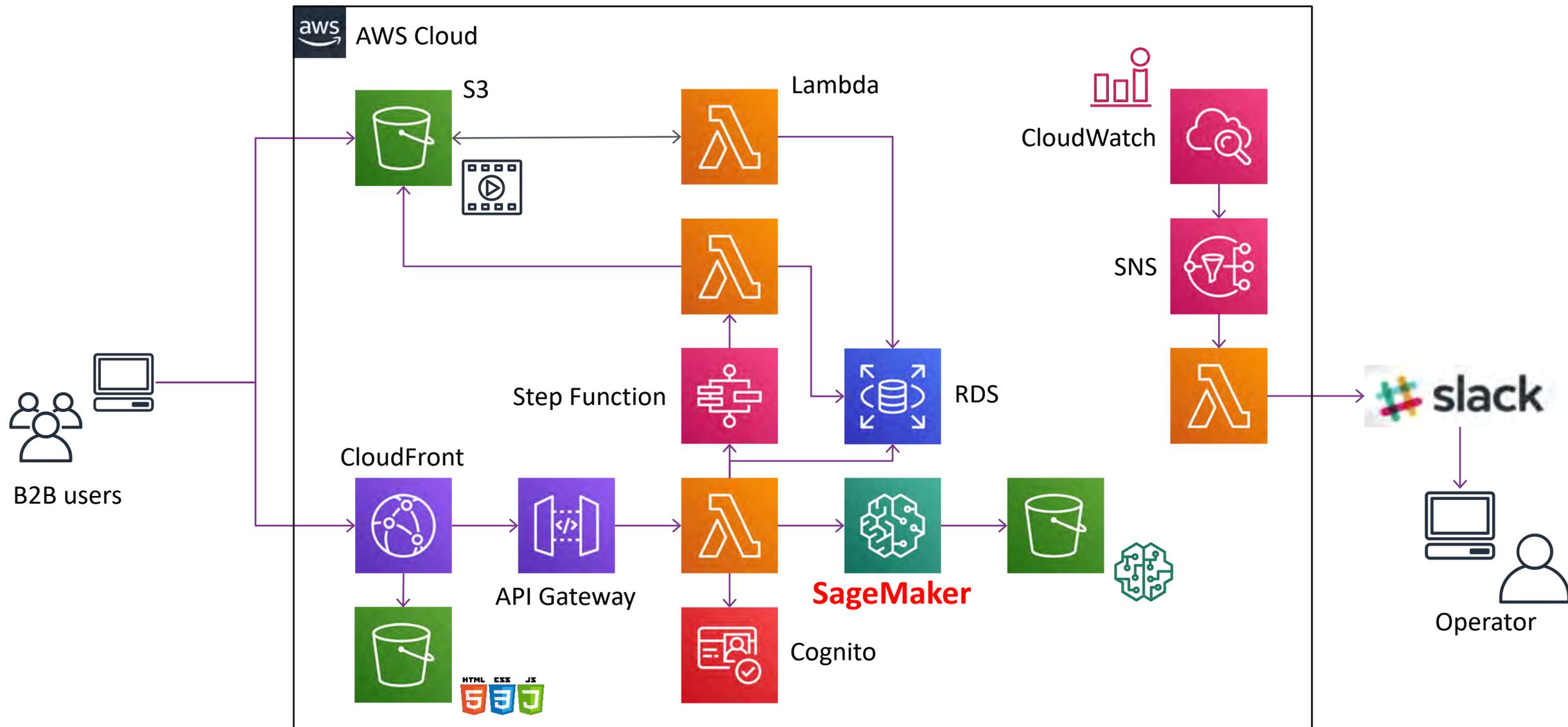


Original

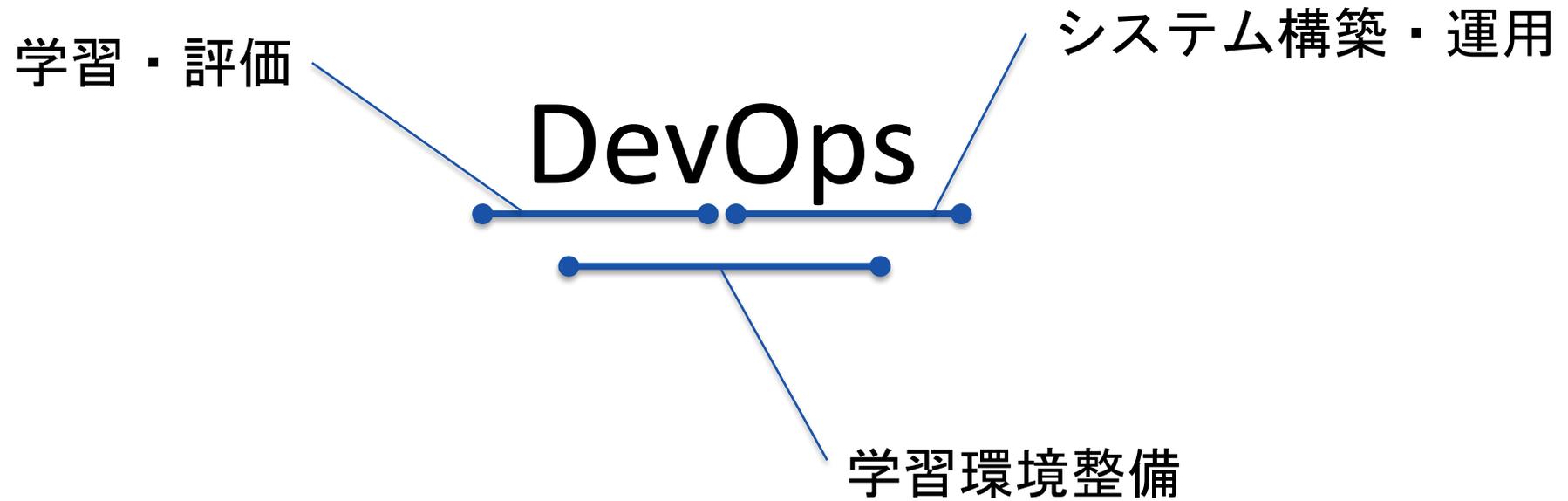


Super Resolution

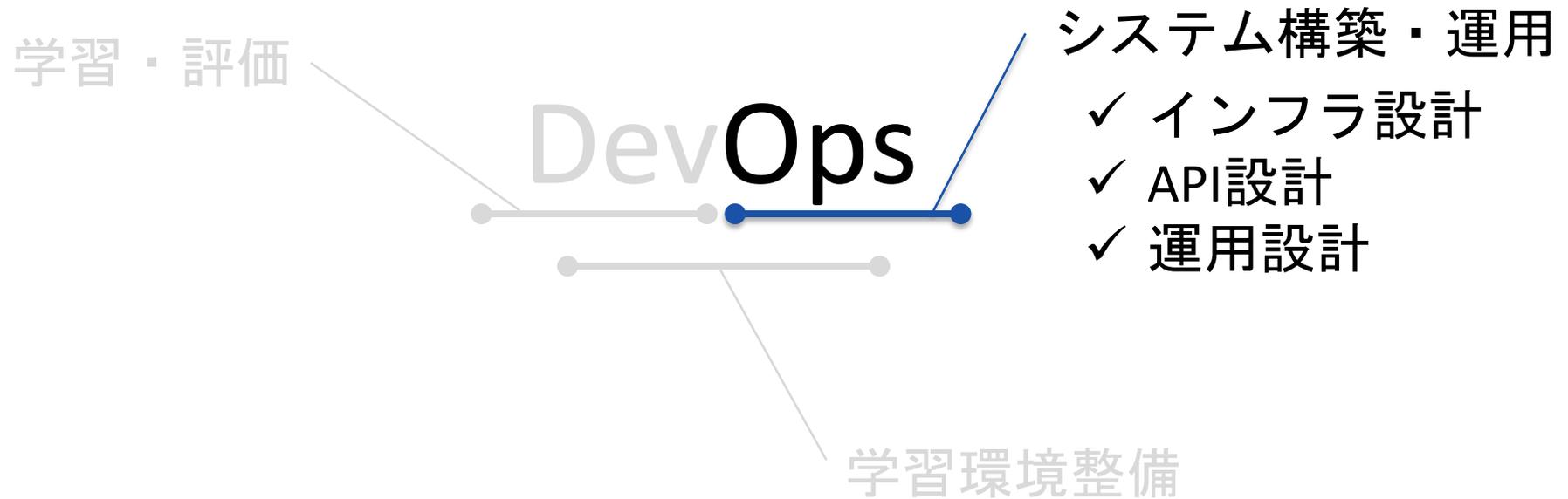
SoVeC Smart Video: システム概要



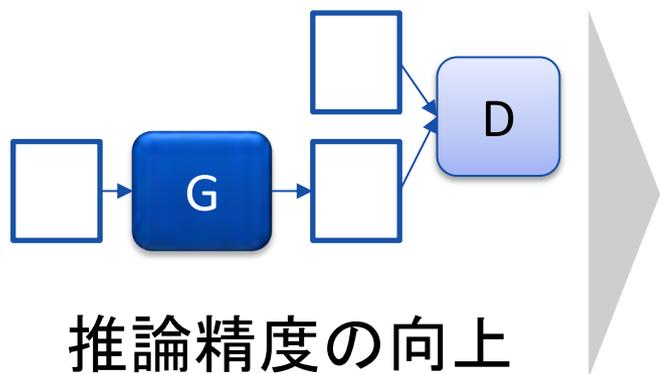
機械学習のDevOps



機械学習のDevOps: システム構築・運用



機能提供方法も課題



機能の提供

インフラ設計 / Computing resource案



ECS

OSレイヤーの運用必要



Fargate

GPUインスタンスなし



Lambda

スペック不足

(2020.02.27 現在)

- ✓ インフラ構成どうする？
- ✓ どのようにI/Fを設計する？
- ✓ 運用どうする？

実装: cf. Edge(前回): 同じリソース上に配置

商用向け運用設計



SageMaker / endpoint機能

```
In [3]: from sagemaker.pytorch.model import PyTorchModel
        pytorch_model = PyTorchModel(
            model_data='s3://sagemaker-us-west-2-1639614456xx/sagemaker-pytorch-2019-10-09-05-11-06-629/outp
            role=role,
            framework_version='1.1.0',
            entry_point='train.py',
            source_dir='./src')

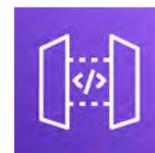
In [ ]: predictor = pytorch_model.deploy(
        initial_instance_count=1,
        instance_type='ml.p2.xlarge')
```

```
predictor = pytorch_model.deploy(
    initial_instance_count=1,
    instance_type='ml.p2.xlarge'
)
```

ミニマムのシステム構築例

- ✓ サーバーレス/フルマネージド
- ✓ 充実したSDKやサンプル
- ✓ シンプルなデプロイ

シンプルなインフラ構成案



API G/W



Lambda



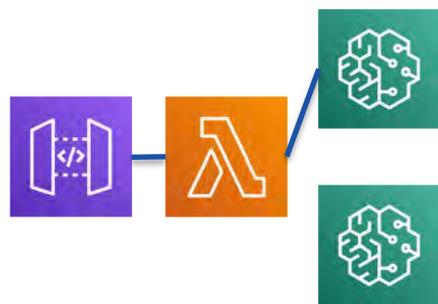
SageMaker

具体的な標準メソッド

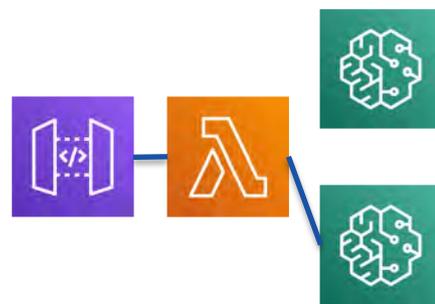
```
def model_fn(model_dir): # モデル読み込み
def input_fn(in, content_type): # 入力データ受取り
def predict_fn(data, model): # 推論の実行
def output_fn(prediction, accept): # レスポンス生成
```

商用運用に耐えられるデプロイ (次頁)

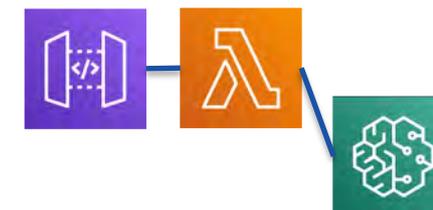
Blue/Greenデプロイを実現



新たなendpointを生成



endpointを変更

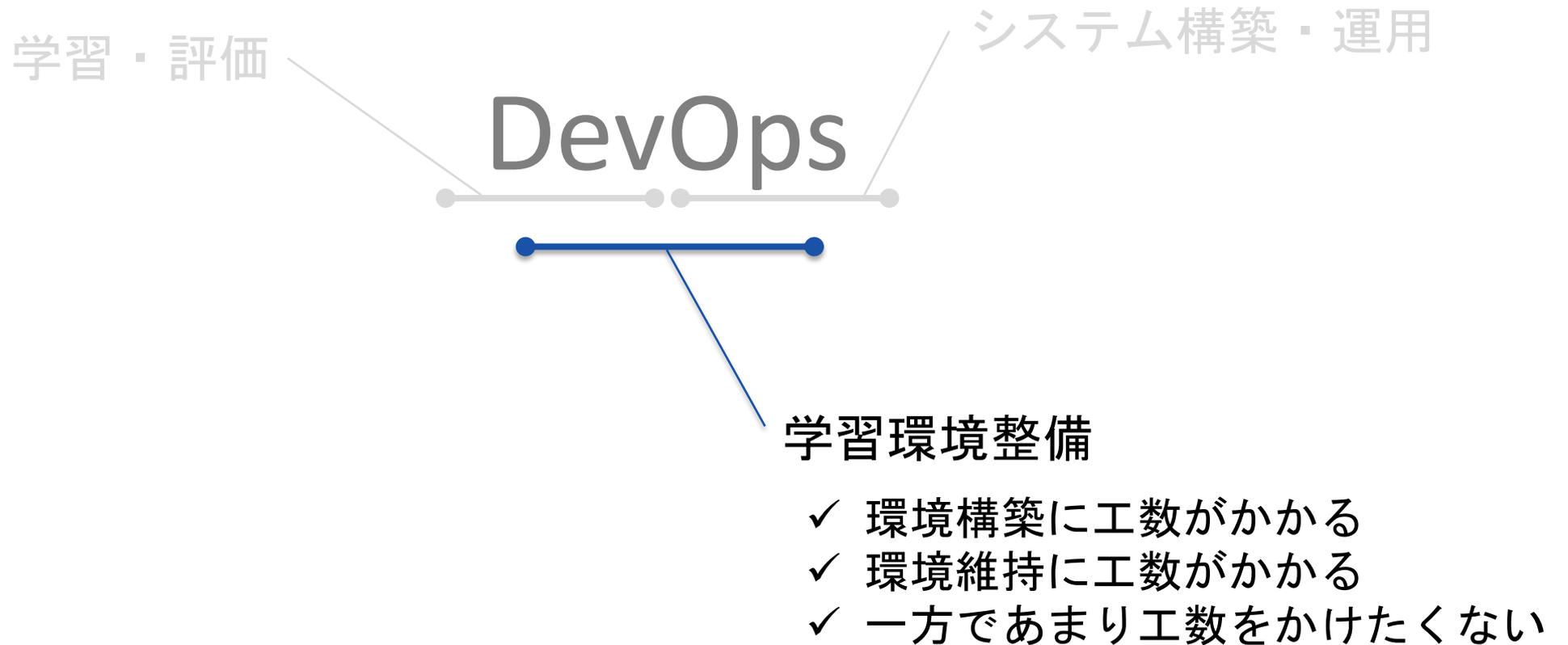


古いendpointを削除

ベストプラクティスなのか？

- ✓ マルチモデルエンドポイント...は違う気がする
- ✓ より良い運用方法があるのではないか？

機械学習のDevOps: 学習環境整備

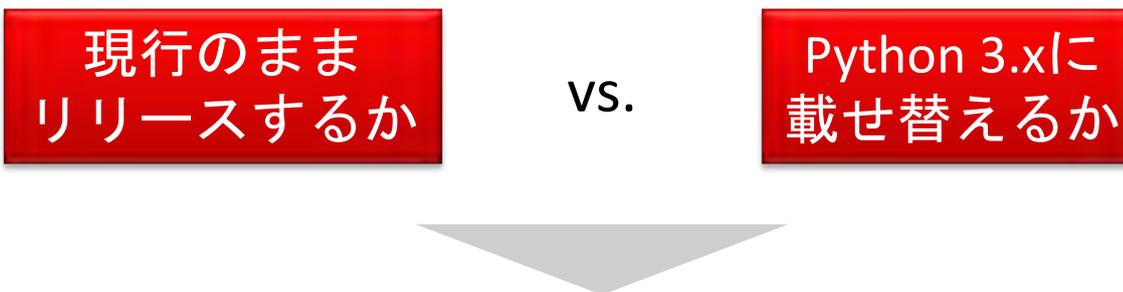




Python 2.7サポート終了のお知らせ / 2019.09.10 (立上げ終盤)

- ✓ 2020年1月1日、Python 2.7 サポート終了
- ✓ Deep Learningの環境を Python 2.7で構築済み
- ✓ 順次サポートが打ち切られるはず...

突きつけられる選択肢



ユーザのデータが入ってきた後は、更に載せ替えに手間がかかるので初回リリース前の載せ替えを選択

学習環境整備: Before

DevOps



単純なPython 2.7 → 3.x だけ変更に残まらない影響範囲（環境構築の煩雑さ）

- ✓ Ubuntu 18.04は意外とHWと相性悪い
- ✓ GPUの進化に対するSWの追従
- ✓ CUDAの最新版は非対応のDL F/W
- ✓ cuDNNとCUDAのバージョンの関係
- ✓ Anacondaとpipの混合
- ✓ DL Framework: 進化が早い
- ✓ 周辺モジュール: 細かく動作しない
- ✓ etc. etc...

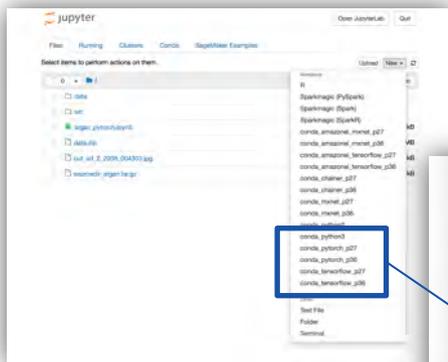
環境の構築・維持に工数が取られる



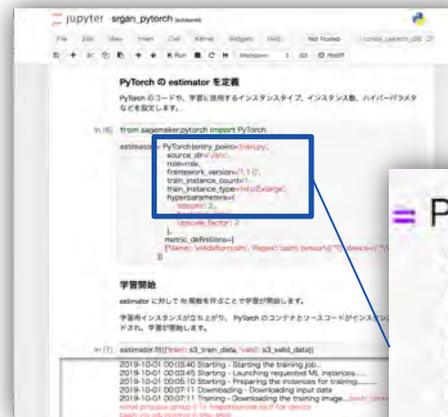
SageMaker / Notebook



SageMaker / Jobs



```
conda_python3
conda_pytorch_p27
conda_pytorch_p36
conda_tensorflow_p27
conda_tensorflow_p36
```



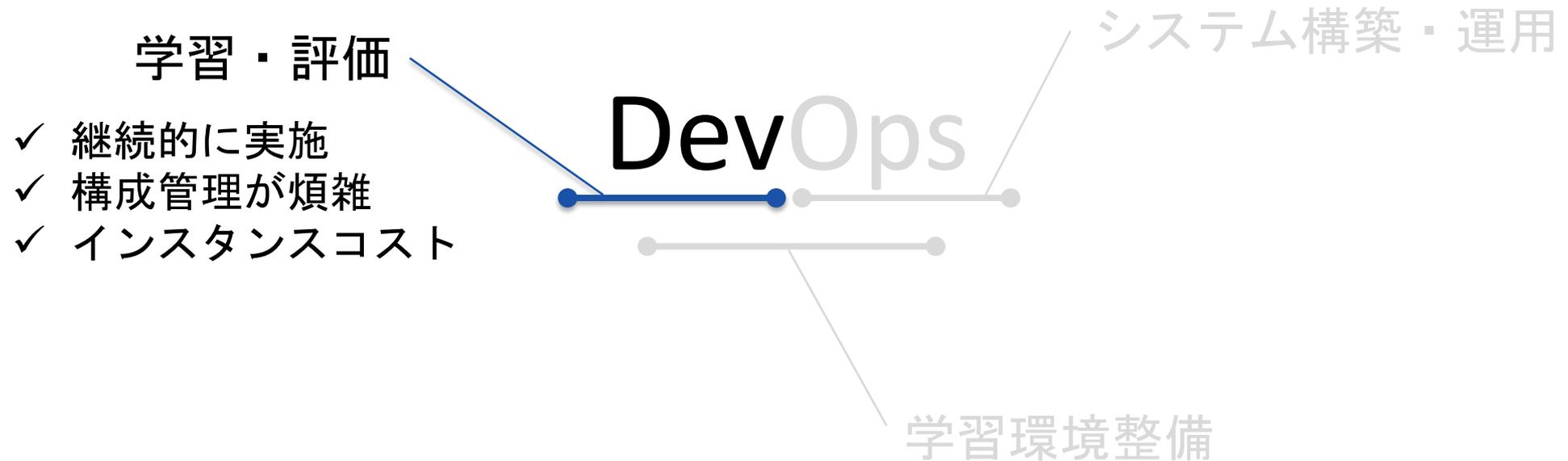
```
= PyTorch(entry_point='train.py',
source_dir='./src',
role=role,
framework_version='1.1.0',
train_instance_count=1,
train_instance_type='ml.p2.xlarge',
hyperparameters={
'epochs': 2
```

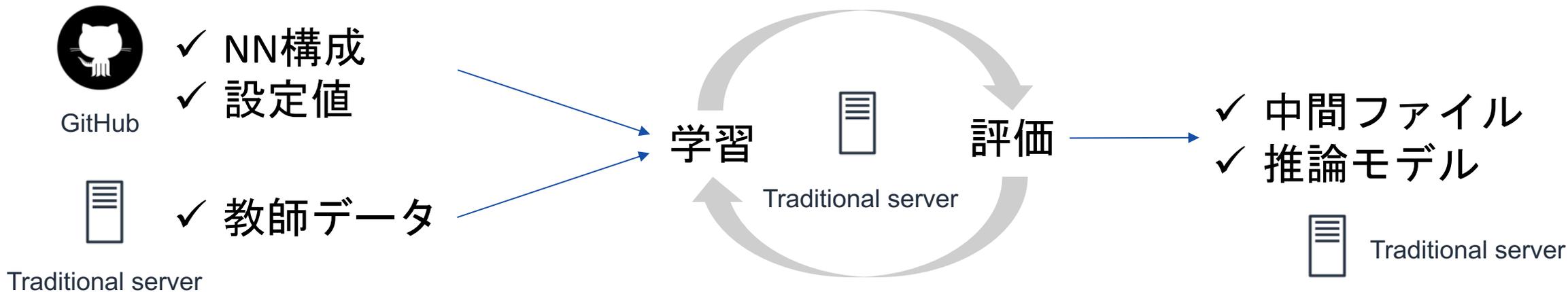
- ✓ ベースのSW stackを選択
- ✓ DL Frameworkのバージョンを選択
- ✓ (もちろん詳細な設定も可能)



学習・評価に
注力できる

機械学習のDevOps: 学習・評価





機械学習の精度向上の難しさ

- ✓ 勾配消失問題
- ✓ 学習効率のためのデータ整備
- ✓ etc. etc....

試行錯誤の連続

構成管理が煩雑

- ✓ 中間ファイルと学習過程の紐付け
- ✓ 前処理した教師データ群の整理
- ✓ 推論モデルとリリース物の紐付け



GitHub

- ✓ NN構成
- ✓ 設定値



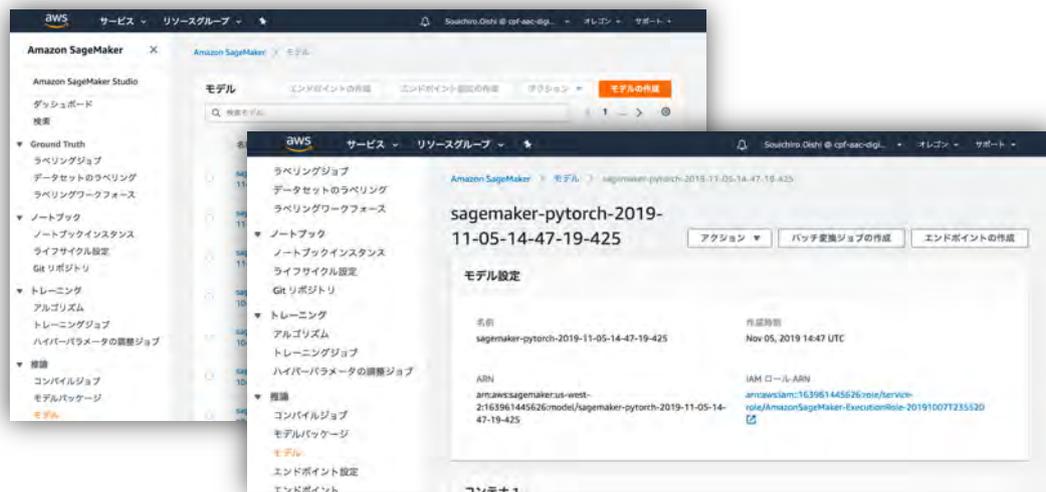
- ✓ 教師データ

学習



評価

- ✓ 中間ファイル
- ✓ 推論モデル



構成管理が容易に

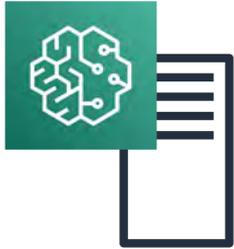
- ✓ in/outの全データをURIで紐付け
- ✓ 教師データのバージョン管理も可能
- ✓ 再デプロイ、ロールバックも容易

機械学習のDevOps: Tips、リクエスト

DevOps

Tips、リクエスト





Local mode: ローカルPCにSageMakerの環境が構築できる

ローカル環境の構築・維持が必要だが下記ユースケースでは有効

- ✓ 会社の規定などの理由で教師データが社外環境に置けない
- ✓ 学習環境が既に整備されている（ホスティングとの連携）
- ✓ 個人ユース

```
In [5]: from sagemaker.tensorflow import TensorFlow

mnist_estimator = TensorFlow(entry_point='mnist.py',
                             role=role,
                             framework_version='1.12.0',
                             training_steps=100,
                             evaluation_steps=100,
                             train_instance_count=2,
                             train_instance_type='local')

mnist_estimator.fit(inputs)
```

instance_type を 'local' とする

Memo: HyperParameter最適化は未対応

```
~/anaconda3/envs/py36sagelocal/lib/python3.6/site-packages/sagemaker/session.py in compile_model(self, input_model_config, output_model_config, role, job_name, stop_condition, tags)
    427
    428     LOGGER.info("Creating compilation-job with name: %s", job_name)
-> 429     self.sagemaker_client.create_compilation_job(**compilation_job_request)
    430
    431     def tune( # noqa: C901

AttributeError: 'LocalSagemakerClient' object has no attribute 'create_compilation_job'
```

(2020.02.27 現在)



Amazon SageMaker で、機械学習のトレーニングコストを最大 90% 削減するマネージドスポットトレーニングが利用できるようになりました

投稿日: Aug 26, 2019

Amazon SageMaker において、マネージドスポットトレーニングと呼ばれる、新しいフルマネージドオプションがサポートされました。これは、Amazon EC2 スポットインスタンスを利用して機械学習モデルをトレーニングするものです。スポットインスタンスでは、AWS クラウド内の使用されていないコンピューティング性能を活用できます。結果として、機械学習モデルのトレーニングコストを最適化し、オンデマンドインスタンスに比べて最大 90% 削減できます。

スポットインスタンスは Amazon SageMaker が管理するため、お客様がキャパシティーを継続的にポーリングする必要はありません。また、追加でツールを構築する必要もありません。トレーニングジョブは、スポットキャパシティーが入手できたときに Amazon SageMaker によってできる限り高い信頼性で実行されます。マネージドスポットトレーニングは、トレーニングモデルが SageMaker 内の一般的な ML フレームワーク、SageMaker の組み込みアルゴリズム、カスタム構築モデルを使用して構築されている場合に使用できます。また、マネージドスポットトレーニングの自動モデルチューニングを使用して、機械学習モデルを調整することもできます。

マネージドスポットトレーニングは、Amazon SageMaker で使用できるすべてのインスタンスタイプでサポートされ、Amazon SageMaker が提供されている AWS リージョンすべてで使用できます。詳細については、AWS ドキュメントを参照してください。また、ブログ記事でも詳細をご覧ください。

<https://aws.amazon.com/jp/about-aws/whats-new/2019/08/amazon-sagemaker-launches-managed-spot-training-saving-machine-learning-training-costs/>

```
estimator = PyTorch(entry_point='train.py',
                    source_dir='./src',
                    role=role,
                    framework_version='1.1.0',
                    train_instance_type='ml.p2.xlarge',
                    train_instance_count=1,
                    train_use_spot_instances=True,
                    train_max_run=36000,
                    train_max_wait=36000,
                    hyperparameters={
                        'epochs': 10000,
                        'backend': 'gloo',
                        'upscale_factor': 2
                    },
                    metric_definitions=[
                        {'Name': 'validation_loss'}
                    ])

train_instance_count=1,
train_use_spot_instances=True,
train_max_run=36000,
train_max_wait=36000,
```

試してみたが...

- ✓ 時間がないときに、なければ待つ、はできない
- ✓ 時間がないときに、中断するかも、は辛かった
- ✓ 時間に余裕があればコストインパクト: 大



Amazon Elastic Interface

- ✓ CPUインスタンスにGPUを付与できる
- ✓ GPUインスタンスよりもコストダウン可能
- ✓ TensorFlowとApache MXNet、ONNXをサポート

(2020.02.27 現在)

Open Neural Network Exchange (ONNX) フォーマットのサポート

ONNX は、深層学習フレームワークでモデルをトレーニングし、推論のために別のモデルに転送できるようにするオープンフォーマットです。これにより、さまざまなフレームワークの相対的な強みを活用することができます。例えば、ONNX を使用すると、PyTorch の柔軟性によって構築およびトレーニングされたモデルを Apache MXNet に転送して、大規模な推論を効率的に実行できるようになります。ONNX は PyTorch、MXNet、Chainer、Caffe2、Microsoft Cognitive Toolkit に統合されており、TensorFlow を含む他の多くのフレームワ

是非PyTorchもサポートしていただきたいです！

<https://aws.amazon.com/jp/machine-learning/elastic-inference/features/>

まとめ

手間のかかる機械学習の環境構築や構成管理などを少ない工数で実現することができ、学習・評価の作業に注力できた。さらに、デプロイ、トレーニング、ビルドといった単位でSageMakerの各機能を段階的に導入できることが良かった。

SONY

SONYはソニー株式会社の登録商標または商標です。

各ソニー製品の商品名・サービス名はソニー株式会社またはグループ各社の登録商標または商標です。その他の製品および会社名は、各社の商号、登録商標または商標です。