

Intel® OpenVINO[™] ツールキットによる AIハンズオンセミナー ~ディープラーニングモデルの簡易利用と推論の高速化~

インテル株式会社

APJデータセンター・グループ・セールス

AIテクニカル・ソリューション・スペシャリスト / AI&アナリティクス・アドボケイト

大内山 浩

@特別ハンズオン for 大阪大学様 (2020/7/21)

本日のスケジュール

13:00 - 13:05	開会
13:05 - 13:35	IntelのAI製品概要およびOpenVINO ツールキットのご紹介
13:35 - 13:50	ハンズオンの説明と準備
13:50 - 14:50	OpenVINOハンズオンPart1 事前学習済みモデルを用いた推論実行
14:35 - 14:50	休憩
14:50 - 16:20	OpenVINOハンズオンPart2 カスタムモデルを用いた推論実行およびモデルの量子化
16:20 - 16:25	閉会



Intel Al Overview

INTELのAI製品概要、および、 OPENVINO ツールキットのご紹介





AI "も" 動かすCPU

あらゆるワークロードに対応できる<u>汎用性と柔軟性</u>がCPUの特徴です。

インテル° など... SSAS. SAP Microsoft IBM ORACLE MARAZON CO Google Cloud TERADATA & Kubeflow Sport O Color Domino DL Studio 統合ワークフロー ディープラーニング 収集、統合、ETL、ELT メタデータの管理 推論の導入 オープン オープン(管理対象) 自社開発 Waterline Data 🏥 collibra **F**TensorFlow **mxnet** Caffe **OpenVINO** など… & kafka TIBCO Blue River ZAL©NI soork BIGDL सं PaddlePaddle Jaspersoft など... 視覚化 pentaho INFORMATICA PYTÖRCH 🕼 ONNX 🛛 také... Spark Streaming talend など... データの前処理 IP(y): IPython Deteractive Computing ggplot2 mold3 マシンラーニングとアナリティクス 🙄 DataRobot Alation 🔊 Datameer ビッグデータの保存と管理 📕 kibana 🔌 matplotlib Lavastorm オープン (管理対象) D DATAWATCH New Paxata pandas Spork MLib オープン 🍊 Boket Rimahout cloudera MAPR 🚳 unifi DataKitchen 📿 TRIFACTA Gephi 🧔 Grafana Flink HORTONWORKS H_O.ai presto 🔅 XGBoost NumPy mongoDB Qu bole 🏷 panoply ClearStory 3 Data-Driven Documents - Burns tomr alteryx 3 STORM 自社開発 ↓↓↓ 自社開発 🚼 Lore IO composable など... 自社開発: 触 Paxata など... MATLAB Cognitive Scale feedzai TIBC 🛟 + a b | e a u Amenity Analytics IT システム管理 Spotfire など... avaamo V Vagrant 🕋 CHEF 📓 STRATOSCALE sentient katacontainers FIS alteryx QPalantir puppet API CLOUD FOUNDRY 🔞 MESOS Jira Software **XEN** TASCIENCE.COM ANSIBLE bluedata など... **≫** MIRANTIS (intel) intel? S s GE SEMANTICS など... New Relic. Jenkins エンタープライズ・アプリケーション **XEON** CORE 17 PLATINUM inside' inside"





ディープラーニング高速化の要 ~AVX-512 & DL Boost~

AVX-512(SIMD)が搭載され、並列演算性能の向上に寄与しております。更に、Deep Learning BoostというAI専用命令により更なるアクセラレーションが期待できます。

Intel® AVX-512

はいってる

(Intel® Advanced Vector Extensions 512)

はいってる



+

Intel® DL Boost

(Intel® Deep Learning Boost)



Skylake世代からAVX-512搭載 Cascade Lake世代からDL Boost搭載



第10世代Ice Lakeから搭載

インテル® AI ソフトウェア: マシンラーニングとディープラーニング

D I TO I	Machine learning	Deep LEarning	Management Tools	
Developer IUDIS App Developers SW Platform Developer	ZOO	OpenVINO [®]	Containers	
Topologies & Models Data Scientist	 Intel Distribution for Python Python (SKlearn, Pandas) 	TensorFlow	kubernetes	
Frameworks Data Scientist		OPylorch	Deep Learning Reference Stack	Architect & DevOps
Graph ML Performance Engineer	 Intel Data Analytics Acceleration Library (Intel DAAL) Intel Math Kernel Library (Intel MKL) 	 Intel Machine Learning Scaling Library (Intel MLSL) Intel® Deep Neural Network Library (DNNL) 	Data Analytics Reference Stack	
Kernel ML Performance Engineer	CPU	CPU = 🗆 GPU = [□ FPGA ■ □	専用

Red font products are the most broadly applicable SW products for AI users



Intel® Distribution for Python*

インテルが実装、かつ、最適化した Python、および、周辺ライブラリ

- Numpy
- Pandas
- Scipy
- Scikit-learn
- XGBoost
- TensorFlow
- etc..



Performance results are based on testing as of Judy 9, 2018 and may not reflect all publicly available security updates. See configuration disclosure for details. No product can be abcolutely secure. Software and workloads used in performance tests may have been optimized for performance only on intel interpropressors. Performance tests, such as SySTamsk and Mobilewink are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more complete information, see <u>Performance Benchmark Test Discoure</u>.

Testing by Intel as of July 2018; Configuration: Stock Pythore python 3.65 compiled from sources obtained at pythonore, many 1.1.43, sctip 1.1.0, sctipl-team 0.191 Installed from pipe, Intel Pythore Intel Distribution for Pythore 2019 Gold; python 3.65 intel, 111, nampy 1.1.43, intel, pp36 5, mild 2019.0 intel_103, mild fit 1.0.2 lintle, np114py36 6, mild 2019.0 intel_103, mild fit 1.0.2 lintle, np114py36 6, sctip 1.1.0 intel_np114py36 6, sctip 1.1.0 intel_np114py36 6, sctipl-1.1.0 intel_np114py36 6, sctipl-1.10 intel_np114py36 7, sctipl-1.10 intel_np114py36 3.5. WS configuration for C 5 10x, Hardware: Intel® Vacom Patitum 8124K CPU (#) 3006Hz (2 sockets, 18 cores/socket, HT2); Virtualization full XVM, two 18 corecl?vs available from ocPU. For c4 8x, Hardware: Intel® Xeont|9 CPU E5-2666 v3 @ 200GHz (2 sockets, 10 cores/socket, 112); Virtualization: full XVM or 114py36 (2 sockets, 10 cores/socket, 112); Virtualization: full XVM or 114py36 (2 sockets, 10 cores/socket, 112); Virtualization: full XVM or 114py36 (2 sockets, 10 cores/socket, 112); Virtualization: full XVM or 114py36 (2 sockets, 10 cores/socket, 112); Virtualization: full XVM or 114py36 (2 sockets, 10 cores/socket, 112); Virtualization: full XVM or 114py36 (2 sockets, 10 cores/socket, 112); Virtualization: full XVM or 114py36 (2 sockets, 10 cores/socket, 112); Virtualization: full XVM or 114py36 (2 sockets, 10 cores/socket, 112); Virtualization: full XVM or 114py36 (2 sockets, 10 cores/socket, 112); Virtualization: full XVM or 114py36 (2 sockets, 10 cores/socket, 112); Virtualization: full XVM or 114py36 (2 sockets, 10 cores/socket, 112); Virtualization: full XVM or 114py36 (2 sockets, 10 cores/socket, 112); Virtualization: full XVM or 114py36 (2 sockets, 10 cores/socket, 112); Virtualization: full XVM or 114py36 (2 sockets, 10 cores/socket, 112); Virtualization: full

Intel's compliers may or may not optimize to the same degree for non-intel microprocessors for optimizations that are not unique to intel microprocessors. These optimizations include SE27_SE3_and SSE3_Instruction ests and other optimizations in this are not unique to intel microprocessors. These optimizations intel and compliance to the microprocessors consors on transmatcher sectors and other optimizations intel and complications are and complications intel and complications intel and complications intel and complications intel and complications are and complications and complications are and complications and complications are and complic



Performance results are based on testing as of July 9, 2018 and may not reflect all publicly available security updates. See configuration disclosure for details. No product can be absolutely secure. Software and workloads used in performance tests may have been optimized for performance endy on their interprocessors. Performance tests, such as 575 mark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more complete information, see <u>Performance Benchmark Test Subcause</u>.

Intel's compliers may or may not optimize to the same degree for non-intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SS22, SSE3, and SSSE3 instruction ests and other optimizations in this does not guarantee the valiability, unicruiteness of any optimization on microprocessors coressors not manufactured by Intel. Microprocessors-dependent optimizations in this product are intended for use with Intel microprocessors. Creatina optimizations in this product are intended for use with Intel microprocessors. Creatina optimizations in this product are intended for use with Intel microprocessors. Creating optimizations in this product are intended for use with Intel microprocessors. Creating optimizations in this product are intended for use with Intel microprocessors. Creating optimizations in this product are intended for use with Intel microprocessors. Creating optimizations in this product are intended for use with Intender to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice. Notice revision #20110804. For microcomplete informations, see out <u>Optimizations</u> with the complete information of the product User and Reference Guides for more information regarding the specific instruction sets covered by this notice. Notice revision #20110804.

<u> https://software.intel.com/en-us/distribution-for-python/benc</u>



Intel® Distribution for Python* - 導入方法一覧

Build from Source

https://software.intel.com/en-us/distributionfor-python/choose-download/linux

Anaconda

<u>https://software.intel.com/en-</u> <u>us/articles/using-intel-distribution-for-</u> <u>python-with-anaconda</u>

Pip

<u>https://software.intel.com/en-</u> us/articles/installing-the-intel-distributionfor-python-and-intel-performance-librarieswith-pip-and

Docker Image

<u>https://software.intel.com/en-</u> us/articles/docker-images-for-intel-python

Linux Repositories

YUM

<u>https://software.intel.com/en-</u> <u>us/articles/installing-intel-free-libs-and-</u> python-yum-repo

APT

<u>https://software.intel.com/en-</u> <u>us/articles/installing-intel-free-libs-and-</u> python-apt-repo

https://software.intel.com/en-us/distribution-for-python/choose-do



ディープラーニング・フレームワーク(インテルによる最適



更なる推論高速化のためにモデルの最適化と量子化

最適化:不要なOpsの除去、複数のOpsの統合などによりモデルをスマート化

量子化*:モデル内部の数値表現をFP32→FP16/INT8等に変換することでスリム化



2020年5月現在、Intel Deep Learning Boost (VNNI)が搭載された 第2世代インテル® Xeon® スケーラブル・プロセッサー、第10世代 インテル® Core™ プロセッサー(Ice Lakeのみ)にてより効力を発揮する。



参考值 ディープラーニング推論処理ベンチマーク インテル® Xeon® Gold 6254 プロセッサー @ 2.10GHz (18 cores × 1 sockets) Resnet50 推論スループット(FPS) 2020年3月20日に計測 性能比(倍) Input=224x224, BS=1, 1 stream 8.00 同じHWでも性能向上可 6.00 4.00 2.00 OpenVINO™ **OpenVINO**[™] 0.00 ツールキット 最適化後 Intel TensorFlow* Intel TensorFlow* 最適化後

1.15.0 最適化前

INT8 (量子化後)

注)インテル社員による性能確認のための個人的なベンチマーク結果であり、インテルの公式結果ではありません。

FP32 (量子化前)

1.15.0 最適化前

ツールキット

2020R1

12

2020R1

OpenVINO[™] ツールキット

https://software.intel.com/en-us/openvino-toolkit

画像処理とディープラーニング推論のためのライブラリスイートです。3つの特徴をぜひ ご理解ください。



inte

特徴1. AIパーツとしてのOpenVINO

学習済みモデルとサンプルを多数提供しております。アプリケーション導入の迅速化を 促進します。

OpenVINO[™] ツールキットのインテル® ディストリビューション提供の事前学習済みモデル

- 年齢と性別
- 頭の位置
- 人物検出 眼高 / 高角検出
- 人、車、自転車の検出
- ナンバープレート検出:小型および前面
- 車両メタデータ
- 人体姿勢推定
- 行動認識 エンコーダーとデコーダー

- テキストの検出と認識
- 車両検出
- 小売環境
- 步行者検出
- 歩行者と車両の検出
- 横断者の属性認識
- 感情認識
- 特定の人物をさまざまな動画で識別 標準および拡張
- 顔のランドマーク検出
- 視線推定

- 路側物の識別
- 高度な路側識別
- 人の検出と行動認識
- 人の再識別 極小 / 超高速
- 顔の再識別
- ランドマーク回帰
- スマート・クラスルームのユースケース
- 単一画像超解像 (3 モデル)
- インスタンス・セグメンテーション
- など…

<u>バイナリーモデル</u>

- 顔検出バイナリー
- 歩行者検出バイナリー



14

ResNet50 バイナリー

OpenVINO[™]ツールキットに含まれる すぐに利用可能なインテルの学習済みモデルの例(無償)



https://software.intel.com/en-us/openvino-toolkit/documentation/pretrained-models

© 2019 Intel Corporation. 無断での引用、転載を禁じます。

特徴2:モデルコンパイラとしてのOpenVINO



重要な理由:トポロジー変換に基づく抑制により、ハードウェアに適したデータ型に変換することで、パフォーマンスを最大化。

 重要な理由: インターフェイスは、ハードウェ アのタイプに応じた動的読み込みのプラグイン として実装。複数のコードを実装および管理す ることなく、タイプごとに最適なパフォーマン

GPU=グラフィックス・プロセシング・ユニット / インテル® プロセッサー・グラフィックスが統合されたインテル内BGEUビジョン・アスを実現 リルドロン・プロダクト。FPGA バージョンと 8 つの Myriad™ X バージョン

最適化に関する注意事項

OpenCL および OpenCL ロゴは、Apple Inc. の商標であり、Khronos の許諾を得て使用されています。

* その他の社名、製品名などは、一般に各社の表示、商標または登録商標

© 2019 Intel Corporation.無断での引用、転載を禁じます。



intel

特徴3: ヘテロジニアス・オーケストレータ としてのOpenVINO 複数種類のチップが混在するヘテロジニアス環境においても容易に各チップを切り替

え、または、併用可能です。





お使いのPCで始められるOpenVINOでのDL推論



Open/INO

Windows 10 Ubuntu, CentOS, Yocto Mac OS



© 2019 Intel Corporation. 無断での引用、転載を禁じます。

事例:理化学研究所様 胸部疾患の自動診断用モデルの推論性能向上



intel

既存システムがディープラーニングで大変 身 ◎penVINO

既存のインテル・アーキテクチャー搭載システムに Open/INIO™Wールキットを使用して付加価値を追加

JLK INSPECTION ヘルスケア

インテル® NUC を利用して胸部X線画像のAIスクリーニングをレント ゲン・バス内で実現。OpenVINO™ツールキットにより判定時間が4.23 秒から2.81秒に短縮。 医療のサービスレベル向上に貢献。 **●** ZEROFOX インターネット・セキュリティー

ソーシャル・メディア等におけるデジタル・アタックの検出。ディープフェイクやブランドを毀損する情報の検出。OpenVINO™ツールキットの利用により2.3倍の性能向上、50%の遅延低減を達成。



下水管の点検用動画の診断をマニュアルからAI化。インテル® Xeon ス ケーラブル・プロセッサーとOpenVINO™ツールキットを利用して診断 時間を80%削減。正確性も20%向上。



Success Stories https://intel.com/openvino-success-stories





OpenVINOのダウンロード&インストール

- ダウンロード
 - https://software.intel.com/content/www/us/en/develop/tools/openvino-toolkit/choosedownload.html
- インストール
 - Linux版
 - https://docs.openvinotoolkit.org/latest/_docs_install_guides_installing_openvino_linux.html
 - Windows 10版
 - https://docs.openvinotoolkit.org/latest/_docs_install_guides_installing_openvino_windows.html



OpenVINOの簡単なインストール方法

- Docker版
 - 公式レポジトリ
 - https://hub.docker.com/u/openvino
 - 非公式レポジトリ(by大内山)
 - https://hub.docker.com/u/hiouchiy



Preparation

ハンズオンの説明と準備





ハンズオンについて

ハンズオンは2つのパートに分けて実施します。

Part1では、AIパーツとしてのOpenVINOに着目し、OpenVINOが持つ様々な事前 学習済みモデルの使い方を体験いただきます。

Part2では、モデルコンパイラとしてのOpenVINOに着目し、実際にカスタムモ デルを作成し、そのモデルをOpenVINOによってCPU上で高速化する手順を体験 いただきます。

今回のハンズオン環境





© 2019 Intel Corporation. 無断での引用、転載を禁じます。

inte

セカンドモニターが無い方はこのようにウィンドウを配置ください

	30 = 9 p 8 = > 0 % M X 9 8 = 1 = X 0 0 0 8 0 9 4	• × 0 4 4 ** = = = • • • • • • •	N ■ X G + - □ X		
← → C △ ■ portal.azure.	com/#@19c120chedutechac.onmicrosoft.com/resource/subscriptions/98395642-e8d6-4498-8b7c-e48a77	/0e4153/resourceGroups/Co, 🕶 😭 🚳			
= Microsoft Azure	,p - 0 ソース・サービス、ドキュメントの検索 (G+/)				
ホーム > リソース グループ > Cogr	hitiveRG >	🧃 circuit.intel.com 🛛 🗙 📔	🖝 「スーパーシティ」構想 技術アイテ 🗙 🐝 「国家戦略特別区域法の一部を	× 🛦 sample-shinbori Keys and En × +	- 🗆 X
Cognitive Services ク teま (Ctrl+/) ■ 根葉	講師の	← → C 介 portal.azure.com/#@19c120chedutechac.onmicrosoft.com/resource/subscriptions/98395642-e8d6-4498-8b7c-e48a770e41 ☆ Q # ● : Apps _ Intel _ Government _ Smart City _ Partners ● SFDC ▲ Googleドライブ ▲ iCloud ⊗ EventHub			
 アクティビ Po アクセス利 	ロナチロリンノ 雨されます。キーを共有しないでください キーを完美的に再生成することもお勧めし	😑 Microsoft Azure 🔎 Sea	arch resources, services, and docs (G+/)	▶ 16 ¢ @ ? ⊙ ^{19c12}	Och@edu.tech.ac.jp 既定のディレクトリ
 ₹ 7 		Home > Resource groups > TARG >	5		
♪ 問題の影響と離 リソース管理		sample-shinbori Cognitive Services	Keys and Endpoint		×
📣 クイックスタート		Search (Ctrl+/)	🗧 🖸 Regenerate Key1 🖸 Regenerate Key2		
 キーとエントボイント 価格レベル 	+-2	Sverview			
 → ネットワーク ID サブスクリプションによる課金 プロパティ 	エンドボイント https://hiouchiytextanalytics.cognitiveservices.azure.com/ 場所 ① eastus	 Activity log Access control (IAM) Tags 	These keys are used to access your Cognitive Service API. Do n using Azure Key Vault. We also recommend regenerating these call. When regenerating the first key you can use the second k	iot share your keys. Store them securely- for example, a keys regularly. Only one key is necessary to make an API cey for continued access to the service.	
 ● ロック ● ケンプレートのエクスポート ■ 振祝 		Diagnose and solve problems RESOURCE MANAGEMENT	^{Show Keys} KEY1 上 上		D
11 왕은		4 Quick start	KEY 2		
। स. २ म. 💌 <u>०</u> .	🛤 🧑 🚾 🗸 😢 刘 📽 🥝 🏨 👲 噦 🙆 🗽	 Keys and Endpoint Pricing tier 	FNDPOINT ウィンドウ		
		Networking	https://sample-shina	D	1
		🚔 Identity	LOCATION ①		
		Billing By Subscription	westus2	P	1
		Properties			
		Locks			
		Export template	*		
					17:13

2020/06/19

I A III 🤚 🚷 🥥 💁 👘 💵 😂 🌉 🗟

Chat

Teams

iend

H Files

OpenVINO Hands-on Part1

PART1: **事前学習済みモデルを**用いた推論実行



ハンズオンPart1の概要

- OpenVINOに用意されている下記の事前学習済みモデルを動かします。
 - 画像セグメンテーション
 - 自動彩色(Colorization)
 - 手書き日本語認識
- 主な手順は下記の通りです。
 - 1. 事前学習済みモデルのダウンロード(OpenVINOモデル管理サーバーより)
 - 2. (必要に応じて)モデルをOpenVINOのIR形式に変換(最適化)
 - 3. IR形式のモデルをPythonスクリプトで実行



ステップ.0諸準備

- SSHクライアント(RLoginがおすすめ)でOCTOPUSへログイン
- 下記コマンドを実行する
 - source /octfs/apl/Anaconda3/bin/activate
 - source activate openvino

- Anacondaを起動(Python3.6を使用するため)
- 本セミナー用の仮想環境に切り替え
- source /octfs/apl/Intel/openvino2020.2.120/openvino/bin/setupvars.sh
 OpenVINOの各モジュールへのパスを環境変数に追加
- 下記コマンドを実行し、OCTOPUS上の画像表示用のウィンドウ(ImageMagick)を立ち上げておく
 - display &
- WinSCPをOctopusへ接続しておく





ステップ.1画像セグメンテーション



下記の順番でコマンドを実行してください

 $\mathsf{cd}\,\,\tilde{}$

 $cp\ -r\ /octfs/apl/Intel/openvino2020.2.120/openvino/inference_engine/demos/python_demos/instance_segmentation_demo/\ .$

cd instance_segmentation_demo

python /octfs/apl/Intel/openvino2020.2.120/openvino/deployment_tools/tools/model_downloader/downloader.py --name instance-segmentationsecurity-1025

 $wget\ https://raw.githubusercontent.com/hiouchiy/IntelAI_and_Cloud/master/Azure/handson/my_instance_segmentation.py$

wget https://isuta.jp/wp-content/uploads/2014/05/top3.jpg

ImageMagick上で top3.jpgを表示する。

下記のいずれかのコマンドを実行する。

- フロントエンド: python my instance_segmentation.py -m intel/instance-segmentation-security-1025/FP32/instance-segmentationsecurity-1025.xml --label coco_labels.txt --no_keep_aspect_ratio -i top3.jpg --no_show
- 計算ノード:次ページ参照

ImageMagick上で after_top3.jpgを表示する。

適当な画像を用意してモデルにセグメンテーションさせてみてください。



計算ノードで実行するためのスクリプト

#!/bin/bash

#PBS -q LECTURE

#PBS -y 326

#PBS -I elapstim_req=1:00:00

cd \$PBS_O_WORKDIR

source /octfs/apl/Anaconda3/bin/activate

source activate openvino

source /octfs/apl/Intel/openvino2020.2.120/openvino/bin/setupvars.sh

python my_instance_segmentation.py -m intel/instance-segmentation-security-1025/FP32/instancesegmentation-security-1025.xml --label coco_labels.txt --no_keep_aspect_ratio -i top3.jpg --no_show

補足: Model DownloaderとOpen Model Zoo



python \$OPENVINO_INSTALL_DIR/deployment_tools/tools/model_downloader/downloader.py -print_all ですべての事前学習済みモデルの一覧が確認できます。







ステップ.2 自動彩色(Colorization)



下記の順番でコマンドを実行してください

$\mathsf{cd}\,\,\tilde{}$

 $cp\ -r\ /octfs/apl/Intel/openvino2020.2.120/openvino/inference_engine/demos/python_demos/colorization_demo/\ .$

cd colorization_demo

python /octfs/apl/Intel/openvino2020.2.120/openvino/deployment_tools/tools/model_downloader/downloader.py --list models.lst

python /octfs/apl/Intel/openvino2020.2.120/openvino/deployment_tools/tools/model_downloader/converter.py --name colorization-v2

wget https://raw.githubusercontent.com/hiouchiy/IntelAI_and_Cloud/master/Azure/handson/mycolorization.py

wget https://github.com/hiouchiy/IntelAI_and_Cloud/raw/master/Azure/handson/black_white.jpg

ImageMagick上で black_white.jpgを表示する

下記いずれかのコマンドを実行する

- フロントエンド: python mycolorization.py --coeffs public/colorization-v2/colorization-v2.npy -m public/colorization-v2/FP32/colorization-v2.xml -i black_white.jpg
- 計算ノード:次のページを参照

ImageMagick上で colorized_ black_white.jpgを表示する



計算ノードで実行するためのスクリプト

#!/bin/bash

#PBS -q LECTURE

#PBS -y 326

#PBS -I elapstim_req=1:00:00

cd \$PBS_O_WORKDIR

source /octfs/apl/Anaconda3/bin/activate

source activate openvino

source /octfs/apl/Intel/openvino2020.2.120/openvino/bin/setupvars.sh

python mycolorization.py --coeffs public/colorization-v2/colorization-v2.npy -m public/colorization-v2/FP32/colorization-v2.xml -i black_white.jpg



ステップ.3 手書き日本語文字認識



下記の順番でコマンドを実行してください

 $\mathsf{cd}\,\,\tilde{}$

 $cp\ -r\ /octfs/apl/Intel/openvino2020.2.120/openvino/inference_engine/demos/python_demos/handwritten_japanese_recognition_demo/\ .$

cd handwritten_japanese_recognition_demo

python /octfs/apl/Intel/openvino2020.2.120/openvino/deployment_tools/tools/model_downloader/downloader.py --list models.lst

python handwritten_japanese_recognition_demo.py -m intel/handwritten-japanese-recognition-0001/FP32/handwritten-japanese-recognition-0001.xml -d CPU -cl data/kondate_nakayosi_char_list.txt -i data/test.png

wget https://github.com/hiouchiy/IntelAI_and_Cloud/raw/master/Azure/handson/myname.png

ImageMagick上でmyname.pngを表示する

下記いずれかのコマンドを実行する

- フロントエンド: python handwritten_japanese_recognition_demo.py -m intel/handwritten-japanese-recognition-0001/FP32/handwritten-japanese-recognition-0001.xml -d CPU -cl data/kondate_nakayosi_char_list.txt -i myname.png
- 計算ノード:次のスライドを参照

適当な手書き文字を書いてモデルに認識させてみてください。

計算ノードで実行するためのスクリプト

#!/bin/bash

#PBS -q LECTURE

#PBS -y 326

#PBS -I elapstim_req=1:00:00

cd \$PBS_O_WORKDIR

source /octfs/apl/Anaconda3/bin/activate

source activate openvino

source /octfs/apl/Intel/openvino2020.2.120/openvino/bin/setupvars.sh

python handwritten_japanese_recognition_demo.py -m intel/handwritten-japanese-recognition-0001/FP32/handwritten-japanese-recognition-0001.xml -d CPU -cl data/kondate_nakayosi_char_list.txt -i myname.png



おまけ:その他のモデルも動かしてみましょう

OpenVINO samples & demo(公式ドキュメント)

- <u>https://docs.openvinotoolkit.org/latest/_demos_README.html</u>
- https://docs.openvinotoolkit.org/latest/_docs_IE_DG_Samples_Overview.html
- <u>https://github.com/opencv/open_model_zoo</u>

OpenVINO samples & demo(非公式ドキュメント)

- <u>https://github.com/yas-sim</u>
- <u>https://github.com/hiouchiy</u>
- <u>https://openvino.jp/</u>



OpenVINO Hands-on Part2

カスタムモデルを用いた推論実行およびモデルの量子化



ハンズオンPart2の概要

- 独自のモデルを作成し、それをOpenVINOで動かすことでどの程度推論性能が向上 するかを見ていきます。
- 主な手順は下記の通りです。
 - 1. カスタムモデルの作成(Keras + TensorFlowでResNet50を転移学習)
 - 2. TensorFlowで推論(通常およびCPU最適化の2パターン実行)
 - 3. モデルをOpenVINOのIR形式に変換(最適化)し、OpenVINO上で実行
 - 4. モデルをOpenVINOで量子化(INT8化)し、OpenVINO上で実行



ステップ.0諸準備

- SSHクライアント(RLoginがおすすめ)でOCTOPUSへログイン
- 下記コマンドを叩く
 - source /octfs/apl/Anaconda3/bin/activate
 - source activate openvino
 - source /octfs/apl/Intel/openvino2020.2.120/openvino/bin/setupvars.sh



ステップ.1カスタムモデルの準備 37種類の犬猫を分類するモデル(ResNet50)の作成

下記の順番でコマンドを実行してください

mkdir model_opt_qtz

cd model_opt_qtz

cd [^]

© 2020 Intel Corporation

```
wget https://hiouchiy.blob.core.windows.net/share/data.zip
```

mkdir train_data

```
mv data.zip train_data/
```

```
unzip train_data/data.zip -d train_data/
```

wget https://raw.githubusercontent.com/hiouchiy/IntelAI and Cloud/master/Azure/handson/train_model.py

python train_model.py (注: SciPy not foundエラーが出た場合は、'pip install scipy --user'を実行下さい。)

学習が始まったら強制終了(Ctrl+C)してください。これはCPUが本気を出していいない性能です。以下のように環境変数を追加下さい

KMP_AFFINITY=granularity=fine,compact,1,0 KMP_BLOCKTIME=1 KMP_SETTINGS=1 OMP_NUM_THREADS=24 python train_model.py

学習終了後、tf_model/mobilenet.pbというファイルが出来ていることを確認下さい。





TensorFlow



計算ノードで実行するためのスクリプト

#!/bin/bash

#PBS -q LECTURE

#PBS -y 326

#PBS -I elapstim_req=1:00:00

cd \$PBS_O_WORKDIR

source /octfs/apl/Anaconda3/bin/activate

source activate openvino

source /octfs/apl/Intel/openvino2020.2.120/openvino/bin/setupvars.sh

KMP_AFFINITY=granularity=fine,compact,1,0 KMP_BLOCKTIME=1 KMP_SETTINGS=1 OMP_NUM_THREADS=24 python train_model.py



解説:謎の環境変数

4つの環境変数はoneDNNの中で動いているOpenMP用です。



https://software.intel.com/content/www/us/en/develop/articles/maximize-tensorflow-performance-on-cpu-considerations-and-recommender-

© 2020 Intel Corporation

ステップ. 2 TensorFlowで推論実行



下記の順番でコマンドを実行してください



cd model_opt_qtz

cd [^]

 $wget\ https://raw.githubusercontent.com/hiouchiy/IntelAI_and_Cloud/master/Azure/handson/infer_script.py$

python infer_script.py ---input_graph tf_model/resnet50.pb ---num_images 300

これはCPUが本気を出していいない性能です。同じコマンドを下記の環境変数と共に実行してください。

- フロントエンド: KMP_AFFINITY=granularity=fine,compact,1,0 KMP_BLOCKTIME=1 KMP_SETTINGS=1 OMP_NUM_THREADS=24 python infer_script.py --input_graph tf_model/resnet50.pb --num_images 300
- 計算ノード:次のページを参照

推論性能が向上していることを確認してください。



計算ノードで実行するためのスクリプト

#!/bin/bash

#PBS -q LECTURE

#PBS -y 326

#PBS -I elapstim_req=1:00:00

cd \$PBS_O_WORKDIR

source /octfs/apl/Anaconda3/bin/activate

source activate openvino

source /octfs/apl/Intel/openvino2020.2.120/openvino/bin/setupvars.sh

KMP_AFFINITY=granularity=fine,compact,1,0 KMP_BLOCKTIME=1 KMP_SETTINGS=1 OMP_NUM_THREADS=24 python infer_script.py --input_graph tf_model/resnet50.pb --num_images 300



ステップ.3 OpenVINOで推論実行



下記の順番でコマンドを実行してください

OpenVINO[®]

cd ~

cd model_opt_qtz

python /octfs/apl/Intel/openvino2020.2.120/openvino/deployment_tools/model_optimizer/mo.py --input_model=tf_model/resnet50.pb -- input_shape=[1,224,224,3]

次のいずれかのコマンドで推論実行してください。

- フロントエンド: python infer_script.py --input_graph resnet50.xml --num_images 300 --openvino
- 計算ノード:次のページを参照

推論性能が向上していることを確認してください。



計算ノードで実行するためのスクリプト

#!/bin/bash

#PBS -q LECTURE

#PBS -y 326

#PBS -I elapstim_req=1:00:00

cd \$PBS_O_WORKDIR

source /octfs/apl/Anaconda3/bin/activate

source activate openvino

source /octfs/apl/Intel/openvino2020.2.120/openvino/bin/setupvars.sh

python infer_script.py ---input_graph resnet50.xml ---num_images 300 ---openvino

補足: Model Optimizerのモデル最適化技術

Linear Operation Fusing

- BatchNorm and ScaleShift decomposition: BN layers decomposes to Mul->Add->Mul->Add sequence; ScaleShift layers decomposes to Mul->Add sequence.
- Linear operations merge: Merges sequences of Mul and Add operations to the single Mul->Add instance.
- Linear operations fusion: Fuses Mul and Add operations to Convolution or FullybConnected layers.

Grouped Convolutions Fusing

Specific optimization that applies for TensorFlow* topologies. (Xception*)







Caffe Resnet269 block (from Netscope)

https://docs.openvinotoolkit.org/latest/openvino_docs_MO_DG_Deep_Learning_Model_Optimizer_DevGuide.html © 2019 Intel Corporation el⁾ Al

ステップ. 4 OpenVINOで量子化&推論実行



cd

 ${\tt cd\ model_opt_qtz}$

mkdir train_data/quantization

cp train_data/val/*/* train_data/quantization/

wget https://raw.githubusercontent.com/hiouchiy/IntelAI and Cloud/master/Azure/handson/resnet50 int8.json

wget https://raw.githubusercontent.com/hiouchiy/IntelAI_and_Cloud/master/Azure/handson/val.txt

pot -c resnet50_int8json

以下のいずれかのコマンドを実行

- フロントエンド: python infer_script.py --input_graph "results/resnet50_int8_DefaultQuantization/YYYY-MM-DD_HH-MM-SS/optimized/resnet50_int8.xml" --num_images 300 --openvino
- 計算ノード:次のページを参照

推論性能が向上していますか?

現在OCTOPUSに搭載されいているXeonプロセッサーの世代(Skylake)では、1.3~1.4倍程度の向上が期待できます。

次のCascade Lakeでは、Deep Learning Boost(別名VNNI)という、量子化(ここではINT8化)されたモデルに特化した命令セットが搭載されているため、更なる性能向上が期待できます。



© 2020 Intel Corporation





計算ノードで実行するためのスクリプト

#!/bin/bash

#PBS -q LECTURE

#PBS -y 326

#PBS -I elapstim_req=1:00:00

cd \$PBS_O_WORKDIR

source /octfs/apl/Anaconda3/bin/activate

source activate openvino

source /octfs/apl/Intel/openvino2020.2.120/openvino/bin/setupvars.sh

python infer_script.py --input_graph "results/resnet50_int8_DefaultQuantization/YYYY-MM-DD_HH-MM-SS/optimized/resnet50_int8.xml" --num_images 300 --openvino

Post-training Optimization Tool (POT)

・トレーニング後のモデルに対する最適化ツール(再量子化、スパース化(まだ))

・使い方は JSONで書かれたconfigファイルを渡すとその設 定に従ってINT8モデルを作ってくれます (calibration.pyは YAML)

・バリデーションデータセットにCOCOなども指定できるよう になった反面、configファイルの書式、設定項目が増えて ちょっと難しくなりました。

・Default algorithm(簡単)と Accuracy aware algorithm(面倒 だが精度維持)が選べます

・DL workbenchからも 使えますので、そっちを使ったほう が楽ですが、DL workbenchは すべての設定が設定できる わけではないので注意。





性能比較(on OCTOPUS フロントエンド)

TF(Not Optimize) TF(Optimize) OpenVINO OpenVINO(INT8)

画像一枚当たりの推論時間(ms)

Pre + Post Inference



itel 53

性能比較(Cascade Lake 6252(24core x 1skt))

TF(Not Optimize) TF(Optimize) OpenVINO OpenVINO(INT8) ■ Pre + Post ■ Inference

画像一枚当たりの推論時間(ms)





皆さんのお持ちのモデルをOpenVINOで走らせてみましょう。



55

© 2020 Intel Corporation

おまけ:前処理(画像加工等)の高速化ヒント

libjpeg-turboの利用

https://github.com/hiouchiy/libjpeg-turbo-how-to

Optane Persistent Memoryの利用

https://tech.preferred.jp/ja/blog/mn-3-launch/



各種ガイドおよびオンデマンド動画のご紹介

TensorFlow* 量子化ガイド

https://github.com/IntelAl/tools/releases/tag/v1.0.0

PyTorch* 量子化ガイド

- https://pytorch.org/docs/stable/quantization.html
- https://pytorch.org/blog/introduction-to-quantization-on-pytorch/
- https://pytorch.org/tutorials/advanced/static_quantization_tutorial.html
- https://pytorch.org/tutorials/advanced/dynamic_quantization_tutorial.html

OpenVINO[™] ツールキット・モデル最適化ガイド(Model Optimizer の使い方) オンデマンド配信中

https://docs.openvinotoolkit.org/latest/_docs_MO_DG_Deep_Learning_Model_Optimizer_DevGuide.html

OpenVINO[™] ツールキット量子化ガイド (Post-Training Optimization Toolkit の使い方) •--

https://docs.openvinotoolkit.org/latest/_README.html

深層学習 Deep Diveセッション @Intel Data Centric Innovation Day オンデマンド配信中

> OpenVINO[™] ツールキット ハンズオン @Intel AI Park オンデマンド配信中



ハンズオンは以上です





One Intel - アナリティクス & AI ハードウェア製品



イベント/メディアなどのご紹介



公式Connpassグループ https://intel.connpass.com/



公式Youtubeチャンネル https://www.youtube.com/channel/UCymh jXNUpudGFa-vFRB-Yvw



非公式ですが、、、面白サンプル多数公開中 https://github.com/yas-sim https://github.com/hiouchiy



C++でプログラムを書く際の事前設定事項

下記URLに記載されております。

https://docs.openvinotoolkit.org/latest/openvino_docs_IE_DG_inference_engine_intro.html#com mon_workflow_for_using_the_inference_engine_api

事前準備としては、OpenVINOインストールディレクトリ内の下記スクリプトを実行ください。

source /octfs/apl/Intel/openvino2020.2.120/openvino/bin/setupvars.sh

その後、プログラムの記述およびコンパイル(ビルド)が可能となります。

ビルドに関しては下記URLもご参照ください。

https://docs.openvinotoolkit.org/latest/openvino_docs_IE_DG_Integrate_with_customer_applicat ion_new_API.html#build_your_application



