

AlxBoard で OpenVINO™ を使用して YOLOv8 物体検出モデルを高速化する

この記事は、Medium に公開されている「[Accelerating YOLOv8 Object Detection Model on AlxBoard with OpenVINO™](#)」の日本語参考訳です。原文は更新される可能性があります。原文と翻訳文の内容が異なる場合は原文を優先してください。



1. はじめに

AlxBoard で OpenVINO™ ツールキットを使用して YOLOv8 分類モデルのデプロイメントと評価を行ってみましょう。この記事では、YOLOv8 物体検出モデルの高速化に注目します。

最初に、この記事で紹介しているサンプルコードのリポジトリをダウンロードし、OpenVINO™ の推論エンジンを使用した YOLOv8 向けの開発環境をセットアップします。

[GitHub* リポジトリ](#) (英語)

git clone コマンドを使用してコードのリポジトリをクローンします。

2. YOLOv8 物体検出 OpenVINO™ IR モデルをエクスポートする

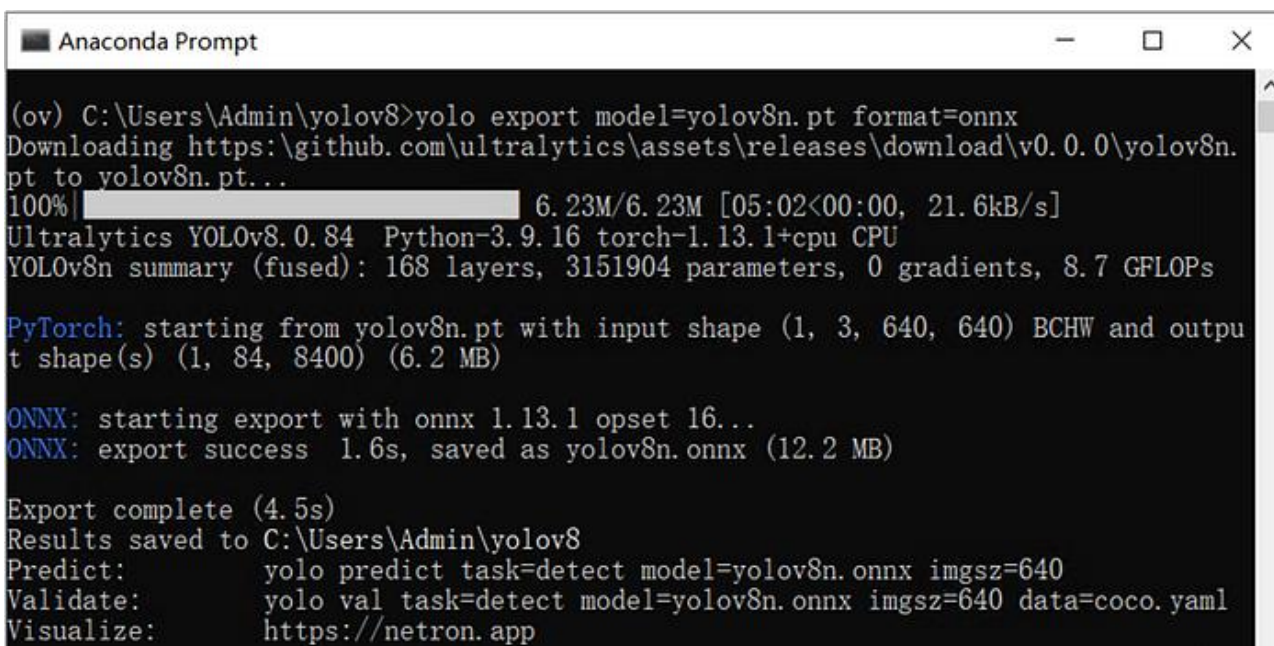
YOLOv8 には、以下の表に示すように、COCO データセットでトレーニングされた 5 つの異なる物体検出モデルがあります。

Model	size (pixels)	mAP ^{val} 50-95	Speed CPU ONNX (ms)	Speed A100 TensorRT (ms)	params (M)	FLOPs (B)
YOLOv8n	640	37.3	80.4	0.99	3.2	8.7
YOLOv8s	640	44.9	128.4	1.20	11.2	28.6
YOLOv8m	640	50.2	234.7	1.83	25.9	78.9
YOLOv8l	640	52.9	375.2	2.39	43.7	165.2
YOLOv8x	640	53.9	479.1	3.53	68.2	257.8

次のコマンドを使用して YOLOv8n.onnx モデルをエクスポートします。

```
yolo export model=yolov8n.pt format=onnx
```

yolov8n.onnx モデルが生成されます。

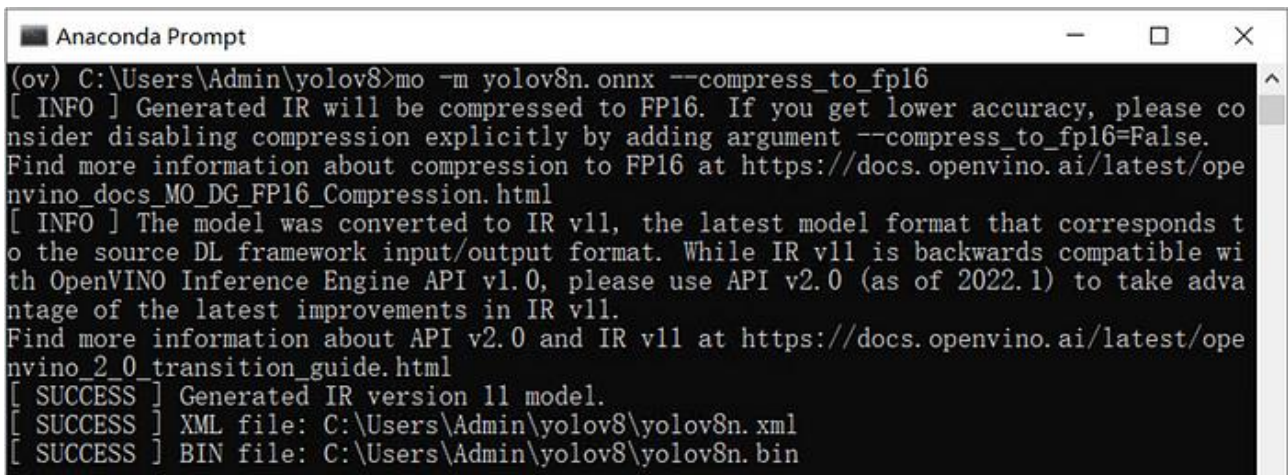


```
Anaconda Prompt
(ov) C:\Users\Admin\yolov8>yolo export model=yolov8n.pt format=onnx
Downloading https://github.com/ultralytics/assets/releases/download/v0.0.0/yolov8n.pt to yolov8n.pt...
100% ██████████ 6.23M/6.23M [05:02<00:00, 21.6kB/s]
Ultralytics YOLOv8.0.84 Python-3.9.16 torch-1.13.1+cpu CPU
YOLOv8n summary (fused): 168 layers, 3151904 parameters, 0 gradients, 8.7 GFLOPs
PyTorch: starting from yolov8n.pt with input shape (1, 3, 640, 640) BCHW and output shape(s) (1, 84, 8400) (6.2 MB)
ONNX: starting export with onnx 1.13.1 opset 16...
ONNX: export success 1.6s, saved as yolov8n.onnx (12.2 MB)

Export complete (4.5s)
Results saved to C:\Users\Admin\yolov8
Predict:      yolo predict task=detect model=yolov8n.onnx imgsz=640
Validate:     yolo val task=detect model=yolov8n.onnx imgsz=640 data=coco.yaml
Visualize:    https://netron.app
```

次のコマンドを使用して、FP16 精度で OpenVINO™ IR 形式モデルを最適化してエクスポートします。

```
mo -m yolov8n.onnx --compress_to_fp16
```



```
Anaconda Prompt
(ov) C:\Users\Admin\yolov8>mo -m yolov8n.onnx --compress_to_fp16
[ INFO ] Generated IR will be compressed to FP16. If you get lower accuracy, please consider disabling compression explicitly by adding argument --compress_to_fp16=False. Find more information about compression to FP16 at https://docs.openvino.ai/latest/openvino_docs_MO_DG_FP16_Compression.html
[ INFO ] The model was converted to IR v11, the latest model format that corresponds to the source DL framework input/output format. While IR v11 is backwards compatible with OpenVINO Inference Engine API v1.0, please use API v2.0 (as of 2022.1) to take advantage of the latest improvements in IR v11. Find more information about API v2.0 and IR v11 at https://docs.openvino.ai/latest/openvino_2_0_transition_guide.html
[ SUCCESS ] Generated IR version 11 model.
[ SUCCESS ] XML file: C:\Users\Admin\yolov8\yolov8n.xml
[ SUCCESS ] BIN file: C:\Users\Admin\yolov8\yolov8n.bin
```

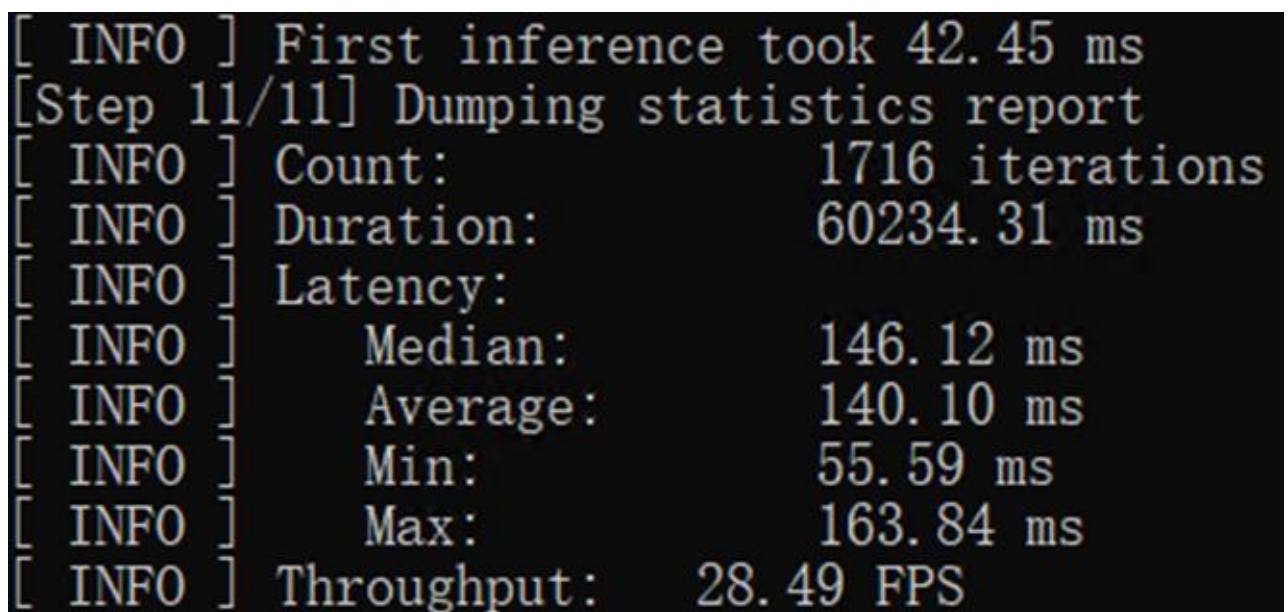
3. benchmark_app を使用して YOLOv8 物体検出モデルの推論パフォーマンスをテストする

benchmark_app は、OpenVINO™ ツールキットで提供されている、AI モデルの推論パフォーマンスを評価するためのパフォーマンス・テスト・ツールです。同期モードまたは非同期モードで、異なる計算デバイスで前処理や後処理を行うことなく、純粋な AI モデルの推論パフォーマンスをテストできます。

次のコマンドを使用します。

```
benchmark_app -m yolov8n.xml -d GPU
```

AlxBoard (英語) の統合 GPU での yolov8n.xml モデルの非同期推論パフォーマンスが表示されます。



```
[ INFO ] First inference took 42.45 ms
[Step 11/11] Dumping statistics report
[ INFO ] Count: 1716 iterations
[ INFO ] Duration: 60234.31 ms
[ INFO ] Latency:
[ INFO ] Median: 146.12 ms
[ INFO ] Average: 140.10 ms
[ INFO ] Min: 55.59 ms
[ INFO ] Max: 163.84 ms
[ INFO ] Throughput: 28.49 FPS
```

4. OpenVINO™ Python* API を使用して YOLOv8 物体検出モデルの推論プログラムを作成する

次の図に示すように、Netron を使用して yolov8n.onnx を開きます。モデルの入力形状は [1,3,640,640]、出力形状は [1,84,8400] です。「84」は、cx、cy、h、w、および 80 クラスのスコアを表します。「8400」は、画像サイズが 640 の場合の YOLOv8 の 3 つの検出ヘッドの出力セル数を示します (80x80+40x40+20x20=8400)。

INPUTS	
images	name: images type: float32[1,3,640,640]

OUTPUTS	
output0	name: output0 type: float32[1,84,8400]

OpenVINO™ Python* API を使用した YOLOv8 物体検出モデルのプログラム例を次に示します。

コアのソースコードは次のとおりです。

```
yolov8_od_ov_sync_infer_demo.py
```

yolov8_od_ov_sync_infer_demo.py の結果を次に示します。



5. まとめ

[AlxBoard](#) (英語) の統合 GPU (24 実行ユニット) と OpenVINO™ を活用することで、YOLOv8 オブジェクト検出モデルのパフォーマンスを高速化できます。[非同期処理](#) (英語) と [AsyncInferQueue](#) (英語) を使用すると、計算デバイスの使用率がさらに向上し、AI 推論プログラムのスループットを向上させることができます。

OpenVINO™ ツールキットとは

AI を加速する無償のツールである OpenVINO™ ツールキットは、インテルが無償で提供しているインテル製の CPU や GPU、VPU、FPGA などのパフォーマンスを最大限に活用して、コンピューター・ビジョン、画像関係をはじめ、自然言語処理や音声処理など、幅広いディープラーニング・モデルで推論を最適化し高速化する推論エンジン/ツールスイートです。

OpenVINO™ ツールキット・ページでは、ツールの概要、利用方法、導入事例、トレーニング、ツール・ダウンロードまでさまざまな情報を提供しています。ぜひ特設サイトにアクセスしてみてください。

<https://www.intel.co.jp/content/www/jp/ja/internet-of-things/opencvino-toolkit.html>

法務上の注意書き

インテルのテクノロジーを使用するには、対応したハードウェア、ソフトウェア、またはサービスの有効化が必要となる場合があります。

絶対的なセキュリティを提供できる製品またはコンポーネントはありません。

実際の費用と結果は異なる場合があります。

© Intel Corporation. Intel、インテル、Intel ロゴ、その他のインテルの名称やロゴは、Intel Corporation またはその子会社の商標です。

* 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です。