

Google Coral Dev Board の分散並列化による AI アクセラレーションの高速化

藤江涼太[†] 並木美太郎[†]

東京農工大学工学部情報工学科[†]

1. はじめに

近年、深層学習の普及に伴いセンサ側で推論を行うエッジAIが注目され始めている。しかし、エッジ側でのAIアクセラレータは消費電力などにより性能が抑えられている。そこで、本研究ではAIアクセラレータの分散並列による推論速度の向上を行い、その台数効果についての検証を行う。

2. 先行研究

先行研究[1]ではAIアクセラレータが搭載されたSoCであるGoogle Coral Dev Board (以下Coral Boardとする)とNvidia Jetson Nanoとの比較評価、およびGPUとの比較を行っている。先行研究によると、Coral BoardはJetson Nanoと比べて電力性能比とメモリ効率が良い。そこで、本研究ではCoral Boardを使用し、分散並列によるさらなる性能向上を目指す。

3. 時間方向への並列による推論速度向上

AIアクセラレータはCNNでの推論が主であり、時系列を扱うRNNへの対応は限定的である。そのため、入力データの分割は、空間方向への分割よりも時間方向の方が容易であり、センサが取得するデータの各フレームを異なるAIアクセラレータで推論するような並列化をおこなえる(図1)。

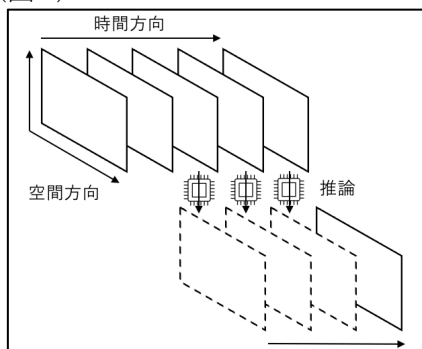


図1 分割方向と推論の並列化

4. Google Coral Dev Board

表1がCoral Boardの主なスペックとなる。この他にSRAM キャッシュとMAC(Multiply Accumulation)のパイプライン化による行列計算への調整と、TensorflowLiteへのサポートを行う。

表1 Coral Board のスペック

CPU	NXP i.MX 8M SOC (quad Cortex-A53, Cortex-M4F)
GPU	Integrated GC7000 Lite Graphics
TPU	Google Edge TPU coprocessor
RAM	1 GB LPDDR4
Flash	8GB eMMC
LAN	Gigabit Ethernet
OS	Mendel Linux 4.0

5. 設計

5.1. 構成

センサノードと推論ノードへ処理を分散し、ネットワークで接続する(図2)。センサノードではセンサデータの配信と推論結果受信を行う。センサノードにはRaspberry Pi 4を、推論ノードにはCoral Boardを用いる。

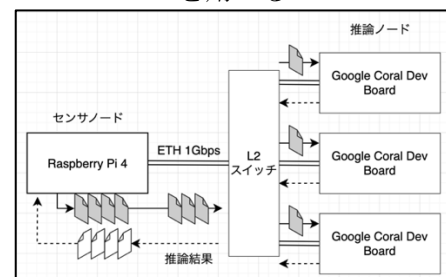


図2 センサノードと推論ノードの接続とデータの配信

5.2. 処理の流れ

T1で入力キューから取得したデータを、T2で推論ノードへTCP/IPで送信し、T3, T4で推論を行った後にT5でセンサノードへ推論結果を送り返し、T6で出力キューへ戻す流れとなる。(図3)

推論ノードの管理にはスレッドを用い、入力キュー・出力キューはスレッド間で共通である。この時の排他処理時間はT1, T6に含まれる。

Improvement AI Acceleration by Distributing and Parallelizing Google Coral Dev Board

[†]Ryota Fujie [†]Mitarou Namiki

[†]Tokyo University of Agriculture and Technology

表2 分散並列数が1台の時の処理時間の内訳

	T1	T2	T3	T4	T5	T6	合計
平均[ms]	0.230	10.712	86.701	23.399	0.069	0.017	121.128
分散[ms]	0.667	1.577	1.463	3.524	0.031	0.004	
割合[%]	0.190%	8.843%	71.578%	19.317%	0.057%	0.014%	

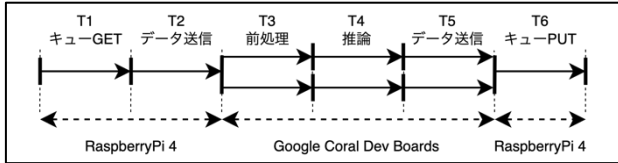


図3 処理の流れ

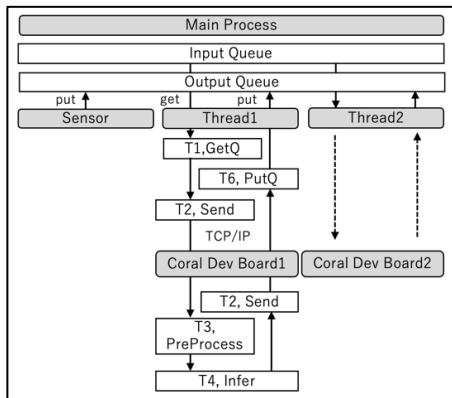


図4 T1からT6の接続

6. 実装

センサノードと推論ノードは全て Python3.9.13 で実装する。Coral Board での推論は Pycoral[2] モジュールを使用し、Raspberry Pi 上でのスレッドとキューには、threading モジュールと queue モジュールを用いる。また、LAN と L2 スイッチは共に 1Gbps である。

なお、推論には 8bit 量子化された MobileNet-V2[3] を使用し、これは物体判定を行うものである。

7. 分散並列化の評価

表2の1台並列時の内訳より、1度の推論処理時間(T3+T4)は0.1101秒であるため、推論処理のみの速度は9.083[FPS]である。また、表3より、オーバーヘッドは、1台並列時で22.93%、2台並列時で23.41%となり、二つの線形近似は式1となる。

表3 2台までの実測値と300秒中のオーバーヘッド

並列数	1	2
FPS	7.00	13.91
オーバーヘッド[sec]	68.79	70.22
割合[%]	22.93	23.41

$$Overhead(n) = 0.004471n + 0.2245 \quad \text{式1}$$

よって、式2より推論処理時間を $InferTime = 0.1101$ として、 n 台並列時の推論速度を推測することができる。

$$\begin{aligned}
 fps(n) &= \frac{n}{InferTime} (1 - 0.2245 - 0.004771n) \\
 &= \frac{n}{0.1101} (0.7755 - 0.004771n) \quad \text{式2} \\
 &= -0.0433317n^2 + 7.04333169n
 \end{aligned}$$

8. 台数効果の評価

表4は、並列数が3台と4台の実測時と、式2より求めた予測値を比較した表である。4台並列時の予測値との相対誤差が0.29%であり、予測に近い実測値を得られることが分かる。これは、式1の定数部(0.2245)を除き、並列化度が99.38%であるからだと考えられる。

また、8台並列であれば、式2より53.57[FPS]となり、7.67倍となることが期待される。

表4 4台までの実測値と式2による推測値

並列数	1	2	3	4
予測値 [FPS]	7.00	13.91	20.74	27.48
実測値 [FPS]	7.00	13.91	20.73	27.56
予測誤差	0.0%	0.0%	0.049%	0.29%

9. おわりに

AI アクセラレータが基本的にサポートするCNNでは時系列を扱わない。そこで、時間方向へ推論の分割と処理の分散により、高い台数効果を得られることを確認でき、そのモデル化を行った。

参考文献

[1] Pilsung KANG, Jongmin JO, Benchmarking Modern Edge Devices for AI Applications, IEICE TRANSACTIONS on Information and Systems, Vol.E104-D, No.3, pp.394-403.
 [2] Google(2023), Pycoral, <https://github.com/google-coral/pycoral>
 [3] Google(2023), Mobilenet-V2 Demo, https://github.com/google-coral/test_data/blob/104342d2d3480b3e66203073dac24f4e2dbb4c41/mobilenet_v2_1.0_224_quant_edgetpu.tflite