

令和元年度 情報工学コース卒業研究報告要旨

枝廣・本田 研究室	氏 名	下平 健太
卒業研究題目	小規模組込みシステム向け DNN フレームワークの比較評価	
<p>近年、多層のニューラルネットワーク(以下 DNN)に代表される深層学習は、技術の発展が進み画像認識や音声認識および自然言語処理などの幅広い分野において使われ注目を集めている。組込みシステムにおいても機器の故障検知などで深層学習の推論の利用が求められている。現状の深層学習を実装するうえで用いられるフレームワーク(以下 DNN フレームワーク)や開発環境は汎用 PC 向けに用意されていることが多く、学習と推論の両方に GPU や FPGA を用いることが一般的である。しかし、小規模な組込みシステムでは消費電力やメモリ使用量などのリソースの制約があるため、多くのリソースが必要なこれらの演算処理装置は使用できず、小規模組込みシステム向けの CPU であるマイコンを使用することになる。マイコンはメモリサイズや性能が低く、例えば最新のマイコンである STM 社の STM32F7 のメモリサイズは最大 2MB で、クロック周波数は 216MHz である。マイコンにおける深層学習の推論を実現するためには C 言語での実装が必要であるため、前述した汎用 PC 向け DNN フレームワークを直接用いて実装することができない。さらに、性能の低いマイコンを用いるためには、メモリ使用量や計算量の削減が可能なパラメータの量子化や SIMD 命令を用いた計算が必要である。以上のような観点から、小規模組込みシステムの DNN の実装は工数がかかるという問題がある。</p> <p>この問題に対して、複数のマイコンベンダーから小規模組込みシステム向け DNN フレームワークがリリースされている。しかしながら、どのフレームワークもリリースから日が浅いため、機能や性能の比較評価が十分なされていない。</p> <p>本研究では、複数の小規模組込みシステム向け DNN フレームワークの機能を調査し比較評価した。具体的には、STM 社が提供している X-CUBE-AI と、NXP 社が提供している eIQ の 2 つについて調査し比較評価した。まず、小規模組込みシステムにおける DNN の実装に必要な観点について整理した。次にそれぞれのフレームワークの特徴や設計フローについて調べ、最後に同一のネットワークに対する性能を比較し考察した。まず機能としてどちらのフレームワークも代表的な汎用 PC 向け DNN フレームワークで作成したネットワークモデルを読み込むことができる。さらに eIQ では、マイコンの持つ SIMD 命令に最適化された DNN ライブラリである CMSIS-NN を用いてネットワーク構造を記述する実装にも対応している。X-CUBE-AI は、量子化を自動で行うことや、推論器のソースコードが自動生成されることなどから実装までの工数が少ないことに対し、eIQ は量子化においてユーザーが別途ツールを用意する必要があることや、最終的なソースコードをユーザーが直接記述しなければならない。</p> <p>同一のネットワークに対する推論処理の性能評価では、ネットワークモデルを読み込ませて実装した場合 X-CUBE-AI は 216MHz のマイコンで 283ms、eIQ は 600MHz のマイコンで 72ms であり、eIQ の CMSIS-NN で実装した場合 96ms だった。また、既存研究では 216MHz のマイコンで同一のネットワークに対する推論を CMSIS-NN で実装した場合は 99ms と報告されており、eIQ の結果はマイコンの性能に対し非常に低速であると分かった。この問題の原因を特定するために、CMSIS-NN とマイコンのアーキテクチャを調査した。eIQ で用いたマイコンは CPU が ROM に対して高速であるため、マイコンとしては大規模な 32Kbytes のキャッシュを搭載している。一方、CMSIS-NN は既存のマイコンの多くがキャッシュを持たないことから SIMD 演算の効率化のみを考慮しており、キャッシュの利用効率の悪いアルゴリズムとなっていたことが原因であった。</p>		