

FPGA 人工知能の ポテンシャルを探る

第3回 FPGA向けニューラル・ネットワークBNNのハードウェア化

ご購入はこちら

鈴木 量三郎

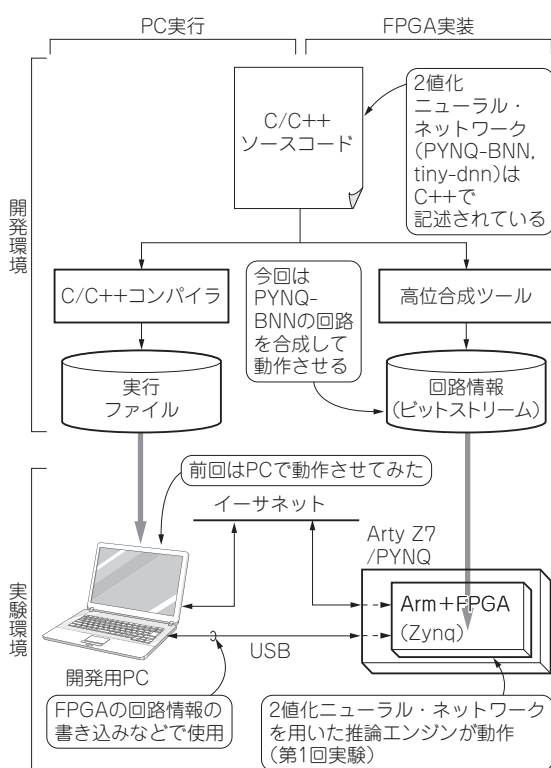


図1 2値化ニューラル・ネットワークBNNの実験環境

今回はPYNQ-BNNを対象に、C++で記述されているコードから回路を合成してボードで動作させてみる

本連載ではFPGAで威力を発揮するといわれている2値化ニューラル・ネットワークBNN (Binarized Neural Networks) のポテンシャルを探ります(図1)。

今回は、サイリンクスが提供しているBNNを詳しく見ていきます。

PYNQ-BNN⁽¹⁾は、サイリンクスのFPGA用に提供されているtiny-dnnを元にした、2値化ニューラル・ネットワークです。

Zynqが内蔵するArmプロセッサ(CPU)からは、メ

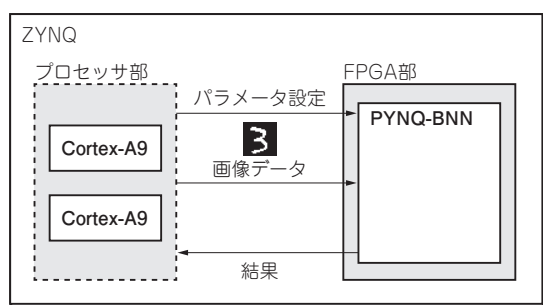


図2 PYNQ-BNNはArmプロセッサ(CPU)からメモリに張り付いたデバイスのように見える

モリに張り付いたデバイスのように見えます(図2)。事前に学習済みのパラメータを設定し、さらに画像を設定することで、結果を得ることができます。

使用するニューラル・ネットワークの構成

● 手書き文字認識MNISTのPYNQ-BNNの構成
lfc-pynqとcnv-pynqの2つのネットワークが用意されています。

MNISTの例で使われているlfc-pynqに注目します(リスト1)。src/network/lfc-pynq/hw/top.cpがFPGAで動作するメイン・ルーチンです。

tiny-dnnの例では、C++の<<演算子をうまく使ってネットワークを表現していました。PYNQ-BNNではネットワークをストリーム化して並列処理させています。ブロックごとに分けてhls::streamというVivado HLSで使う特別な変数でネットワークを構成し、やはりVivado HLSのキーワードであるDATAFLOWを#pragmaで追加しています。

Mem2Stream_BatchとStreamingFCLayer_Batchが処理の中心を担っています。StreamingFCLayer_Batchがニューラル・ネットワークの各層に対応し、C++のtemplate引き数で層の設定をす