

An Adaptable Deep Learning Accelerator Unit (DLAU) for FPGA

N. Sireesha¹ & P.Malleswari²

¹PG Scholar, Dept of ECE, Narsaraopeta Institute of Technology, Yellamanda, Narsaraopeta, Guntur district, Andhra Pradesh, India.

²Assistant Professor, Dept of ECE, Narsaraopeta Institute of Technology, Yellamanda, Narsaraopeta, Guntur district, Andhra Pradesh, India.

Received: May 03, 2018

Accepted: June 11, 2018

ABSTRACT

As the emerging field of machine learning, deep learning shows excellent ability in solving complex learning problems. However, the size of the networks becomes increasingly large scale due to the demands of the practical applications, which poses significant challenge to construct a high performance implementations of deep learning neural networks. In order to improve the performance as well as to maintain the low power cost, in this paper we design deep learning accelerator unit (DLAU), which is a scalable accelerator architecture for large-scale deep learning networks using field-programmable gate array (FPGA) as the hardware prototype. The DLAU accelerator employs three pipelined processing units to improve the throughput and utilizes tile techniques to explore locality for deep learning applications. Experimental results on the state-of-the-art Xilinx FPGA board demonstrate that the DLAU accelerator is able to achieve up to 36.1× speedup comparing to the Intel Core2 processors, with the power consumption at 234 mW.

Keywords: Deep learning, field-programmable gate array (FPGA), hardware accelerator, neural network.

1. INTRODUCTION

In the past few years, machine learning has become pervasive in various research fields and commercial applications, and achieved satisfactory products. The emergence of deep learning speeded up the development of machine learning and artificial intelligence. Consequently, deep learning has become a research hot spot in research organizations [1]. In general, deep learning uses a multilayer neural network model to extract high-level features which are a combination of low-level abstractions to find the distributed data features, in order to solve complex problems in machine learning. Currently, the most widely used neural models of deep learning are deep neural networks (DNNs) [2] and convolution neural networks (CNNs) [3], which have been proved to have excellent capability in solving picture recognition, voice recognition, and other complex machine learning tasks.

However, with the increasing accuracy requirements and complexity for the practical applications, the size of the neural networks becomes explosively large scale, such as the Baidu Brain with 100 billion neuronal connections, and the Google cat-recognizingsystem with one billion neuronal connections. The explosive volume of data makes the data centers quite power consuming. In particular, the electricity consumption of data centers in U.S. are projected to increase to roughly 140 billion kilowatt-hours annually by 2020 [4]. Therefore, it poses significant challenges to implement high performance deep learning networks with low power cost, especially for large-scale deep learning neural network models.

So far, the state-of-the-art means for accelerating deep learning algorithms are field-programmable gate array (FPGA), application specific integrated circuit (ASIC), and graphic processing unit (GPU). Compared with GPU acceleration, hardware accelerators like FPGA and ASIC can achieve at least moderate performance with lower power consumption. However, both FPGA and ASIC have relatively limited computing resources, memory, and I/O bandwidths, therefore it is challenging to develop complex and massive DNNs using hardware accelerators.

For ASIC, it has a longer development cycle and the flexibility is not satisfying. Chen et al. [6] presented a ubiquitous machine-learning hardware accelerator called DianNao, which initiated the field of deep learning processor. It opens a new paradigm to machine learning hardware accelerators focusing on neural networks. But DianNao is not implemented using reconfigurable hardware like FPGA, therefore it cannot adapt to different application demands. Currently, around FPGA acceleration researches, Ly and Chow [5] designed FPGA-based solutions to accelerate the restricted Boltzmann machine (RBM). They created dedicated hardware processing cores which are optimized for the RBM algorithm.

Similarly, Kim et al. [7] also developed an FPGA-based accelerator for the RBM. They use multiple RBM processing modules in parallel, with each module responsible for a relatively small number of nodes. Other similar works also present FPGA-based neural network accelerators [9]. Yu et al. [8] presented an FPGA-based accelerator, but it cannot accommodate changing network size and network topologies. To sum

up, these studies focus on implementing a particular deep learning algorithm efficiently, but how to increase the size of the neural networks with scalable and flexible hardware architecture has not been properly solved.

To tackle these problems, we present a scalable deep learning accelerator unit named DLAU to speed up the kernel computational parts of deep learning algorithms. In particular, we utilize the tile techniques, FIFO buffers, and pipelines to minimize memory transfer operations, and reuse the computing units to implement the largesize neural networks.

This approach distinguishes itself from previous literatures with following contributions.

1) In order to explore the locality of the deep learning application, we employ tile techniques to partition the large scale input data. The DLAU architecture can be configured to operate different sizes of tile data to leverage the tradeoffs between speedup and hardware costs. Consequently, the FPGA-based accelerator is more scalable to accommodate different machine learning applications.

2) The DLAU accelerator is composed of three fully pipelined processing units, including tiled matrix multiplication unit (TMMU), part sum accumulation unit (PSAU), and activation function acceleration unit (AFAU). Different networktopologies such as CNN, DNN, or even emerging neural networks can be composed from these basic modules. Consequently, the scalability of FPGA-based accelerator is higher than ASIC-based accelerator.

Deep learning:

Deep learning (also known as deep structured learning or hierarchical learning) is part of a broader family of machine learning methods based on learning data representations, as opposed to task-specific algorithms. Learning can be supervised, semi-supervised or unsupervised.

Deep learning is a class of machine learning algorithms that:

- Use a cascade of multiple layers of nonlinear processing units for feature extraction and transformation. Each successive layer uses the output from the previous layer as input.
- Learn in supervised (e.g., classification) and/or unsupervised (e.g., pattern analysis) manners.
- Learn multiple levels of representations that correspond to different levels of abstraction; the levels form a hierarchy of concepts.

Deep neural network (DNN):

A deep neural network (DNN) is an artificial neural network (ANN) with multiple hidden layers between the input and output layers. DNNs can model complex non-linear relationships. DNN architectures generate compositional models where the object is expressed as a layered composition of primitives. The extra layers enable composition of features from lower layers, potentially modeling complex data with fewer units than a similarly performing shallow network. The fig.1 shows deep learning neural network.

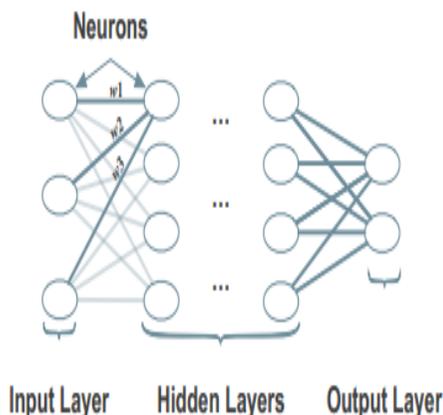


Fig. 1: Deep neural network (DNN) architecture

II. TILE TECHNIQUES AND HOT SPOT PROFILING

RBM's have been widely used to efficiently train each layer of a deep network. Normally, a DNN is composed of one input layer, several hidden layers and one classifier layer. The units in adjacent layers are all-to-all weighted connected. The prediction process contains feedforward computation from given input neurons to the output neurons with the current network configurations. Training process includes pretraining which locally tune the connection weights between the units in adjacent layers and global training which globally tune the connection weights with back propagation (BP) process.

The large-scale DNNs include iterative computations which have few conditional branch operations, therefore, they are suitable for parallel optimization in hardware. In this paper, we first explore the hot spot using the profiler. Results in Fig.2 illustrates the percentage of running time including matrix multiplication (MM), activation, and vector operations. For the representative three key operations: 1) feed forward; 2) RBM; and 3) BP, MM play a significant role of the overall execution. In particular, it takes 98.6%, 98.2%, and 99.1% of the feed forward, RBM, and BP operations. .

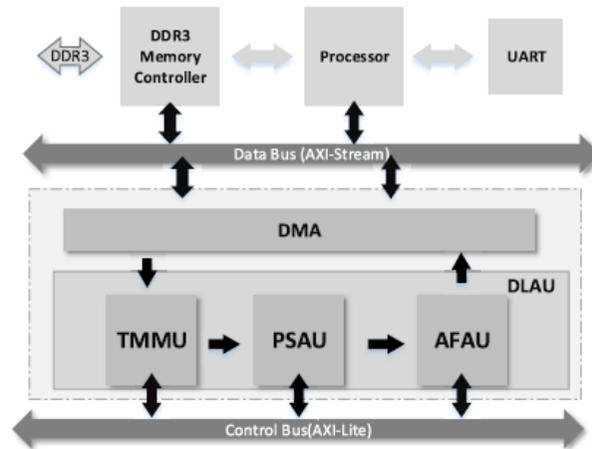


Fig. 2. DLAU accelerator architecture.

In comparison, the activation function only takes 1.40%, 1.48%, and 0.42% of the three operations. Experimental results on profiling demonstrate that the design and implementation of MM accelerators is able to improve the overall speedup of the system significantly. However, considerable memory bandwidth and computing resources are needed to support the parallel processing, consequently it poses a significant challenge to FPGA implementations compared with GPU and CPU optimization measures. In order to tackle the problem, in this paper we employ tile techniques to partition the massive input data set into tiled subsets. Each designed hardware accelerator is able to buffer the tiled subset of data for processing. In order to support the large-scale neural networks, the accelerator architecture are reused. Moreover, the data access for each tiled subset can run in parallel to the computation of the hardware accelerators.

In particular, for each iteration, output neurons are reused as the input neurons in next iteration. To generate the output neurons for each iteration, we need to multiply the input neurons by each column in weights matrix. As illustrated in Algorithm 1, the input data are partitioned into tiles and then multiplied by the corresponding weights. Thereafter the calculated part sum are accumulated to get the result. Besides the input/output neurons, we also divided the weight matrix into tiles corresponding to the tile size.

Algorithm 1 Pseudocode Code of the Tiled Inputs

Require:

```

Ni: the number of the input neurons
No: the number of the output neurons
Tile_Size: the tile size of the input data
batchsize: the batch size of the input data
for n = 0; n < batchsize; n ++ do
  for k = 0; k < Ni; k += Tile_Size do
    for j = 0; j < No; j ++ do
      y[n][j] = 0;
      for i = k; i < k + Tile_Size && i < Ni; i ++ do
        y[n][j] += w[i][j] * x[n][i]
        if i == Ni - 1 then
          y[n][j] = f(y[n][j]);
        end if
      end for
    end for
  end for
end for
end for

```

As a consequence, the hardware cost of the accelerator only depends on the tile size, which saves significant number of hardware resources. The tiled technique is able to solve the problem by implementing large networks with limited hardware. Moreover, the pipelined hardware implementation is another advantage of FPGA technology compared to GPU architecture, which uses massive parallel SIMD architectures to improve the overall performance and throughput. According to the profiling results

depicted in Table I, during the prediction process and the training process in deep learning algorithms, the common but important computational parts are MM and activation functions, consequently in this paper we implement the specialized accelerator to speed up the MM and activation functions.

TABLE I
PROFILING OF HOT SPOTS OF DNN

Algorithms	Matrix Multiplication	Activation	Vector
Feedforward	98.60%	1.40%	
RBM	98.20%	1.48%	0.30%
BP	99.10%	0.42%	0.48%

III. DLAU ARCHITECTURE AND EXECUTION MODEL

Fig.2 describes the DLAU system architecture which contains an embedded processor, a DDR3 memory controller, a DMA module, and the DLAU accelerator. The embedded processor is responsible for providing programming interface to the users and communicating with DLAU via JTAG-UART. In particular it transfers the input data and the weight matrix to internal BRAM blocks, activates the DLAU accelerator, and returns the results to the user after execution.

The DLAU is integrated as a standalone unit which is flexible and adaptive to accommodate different applications with configurations. The DLAU consists of three processing units organized in a pipeline manner: 1) TMMU; 2) PSAU; and 3) AFAU. For execution, DLAU reads the tiled data from the memory by DMA, computes with all the three processing units in turn, and then writes the results back to the memory. In particular, the DLAU accelerator architecture has the following key features.

FIFO Buffer: Each processing unit in DLAU has an input buffer and an output buffer to receive or send the data in FIFO. These buffers are employed to prevent the data loss caused by the inconsistent throughput between each processing unit.

Tiled Techniques: Different machine learning applications may require specific neural network sizes. The tile technique is employed to divide the large volume of data into small tiles that can be cached on chip, therefore the accelerator can be adopted to different neural network size. Consequently, the FPGA-based accelerator is more scalable to accommodate different machine learning applications.

Pipeline Accelerator: We use stream-like data passing mechanism (e.g., AXI-Stream for demonstration) to transfer data between the adjacent processing units, therefore, TMMU, PSAU, and AFAU can compute in streaming-like manner. Of these three computational modules, TMMU is the primary computational unit, which reads the total weights and tiled nodes data through DMA, performs the calculations, and then transfers the intermediate part sum results to PSAU. PSAU collects part sums and performs accumulation. When the accumulation is completed, results will be passed to AFAU. AFAU performs the activation function using piecewise linear interpolation methods. In the rest of this section, we will detail the implementation of these three processing units, respectively.

A. TMMU Architecture

TMMU is in charge of multiplication and accumulation operations. TMMU is specially designed to exploit the data locality of the weights and is responsible for calculating the part sums. TMMU employs an input FIFO buffer which receives the data transferred from DMA and an output FIFO buffer to send part sums to PSAU. Fig. 3 illustrates the TMMU schematic diagram, in which we set tile size = 32 as an example. TMMU first reads the weight matrix data from input buffer into different BRAMs in 32 by the row number of the weight matrix ($n = i\%32$ where n refers to the number of BRAM, and i is the row number of weight matrix). Then, TMMU begins to buffer the tiled node data.

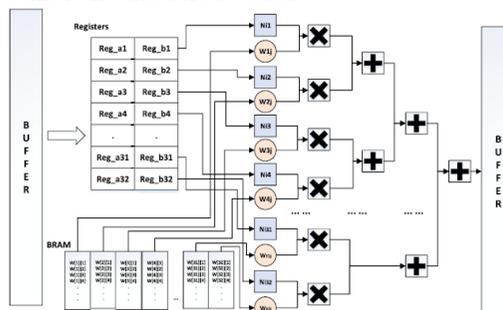


Fig. 3. TMMU schematic.

In the first time, TMMU reads the tiled 32 values to registers Reg_a and starts execution. In parallel to the computation at every cycle, TMMU reads the next node from input buffer and saves to the registers Reg_b. Consequently, the registers Reg_a and Reg_b can be used alternately. For the calculation, we use pipelined binary adder tree structure to optimize the performance. As depicted in Fig. 3, the weight data and the node data are saved in BRAMs and registers. The pipeline takes advantage of time-sharing the coarse-grained accelerators. As a consequence, this implementation enables the TMMU unit to produce a part sum result every clock cycle.

B. PSAU Architecture

PSAU is responsible for the accumulation operation. Fig.4 presents the PSAU architecture, which accumulates the part sum produced by TMMU. If the part sum is the final result, PSAU will write the value to output buffer and send results to AFAU in a pipeline manner. PSAU can accumulate one part sum every clock cycle, therefore the throughput of PSAU accumulation matches the generation of the part sum in TMMU.

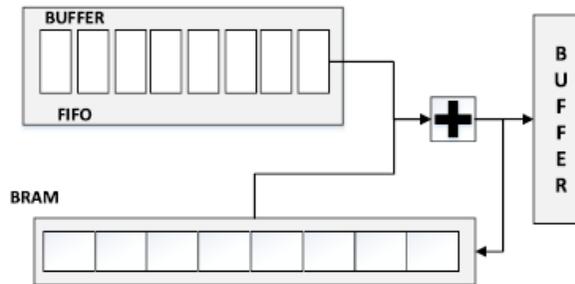


Fig. 4. PSAU schematic.

C. AFAU Architecture

Finally, AFAU implements the activation function using piecewise linear interpolation ($y = a_i * x + b_i$, $x \in [x_i, x_{i+1}]$). This method has been widely applied to implement activation functions with negligible accuracy loss when the interval between x_i and x_{i+1} is insignificant. Equation (1) shows the implementation of sigmoid function. For $x > 8$ and $x \leq -8$, the results are sufficiently close to the bounds of 1 and 0, respectively. For the cases in $-8 < x \leq 0$ and $0 < x \leq 8$, different functions are configured. In total, we divide the sigmoid function into four segments

$$f(x) = \begin{cases} 0 & \text{if } x \leq -8 \\ 1 + a \left[\left[\frac{-x}{k} \right] \right] x - b \left[\left[\frac{-x}{k} \right] \right] & \text{if } -8 < x \leq 0 \\ a \left[\left[\frac{x}{k} \right] \right] x + \left[\left[\frac{x}{k} \right] \right] & \text{if } 0 < x \leq 8 \\ 1 & \text{if } x > 8. \end{cases} \quad (1)$$

Similar to PSAU, AFAU also has both input buffer and output buffer to maintain the throughput with other processing units. In particular, we use two separate BRAMs to store the values of a and b. The computation of AFAU is pipelined to operate sigmoid function every clock cycle. As a consequence, all the three processing units are fully pipelined to ensure the peak throughput of the DLAU accelerator architecture.

IV. IMPLEMENTATION DETAILS

DLAU is designed to accelerate the prediction process of a DNN application. The main computation part is matrix multiplication and activation function. And matrix multiplication is suitable for parallel processing.

A. Time-sharing computation:

For large scale deep learning neural networks, we cannot map all the parallel computation. According to the number of hardware resources, we implement only partial computational logic, and use time-sharing technology to complete the total computation.

1) For several or many layers: As the deep learning neural networks computes layer by layer from input layer to output layer, and the outputs of current layer are the inputs of the next layer, We only implement the largest layer which have the max size of weight matrix. When a smaller layer is in processing, the remaining space in the input or weight coefficients block is padded with zeros.

2) For single layer: We implement partial arithmetic logics. In our case, we design 31 floating point addition and 32 floating point multiplication, then we decompose the input data into data fragment which size is 32. Every iteration the accelerator deals with 32 input values to do the dot product calculation.

B. Data access optimization

Because we need upload the computation task to FPGA, the data access between memory and FPGA is unavoidable. To get better performance, we must optimize the number and time of the data access. We use the BRAM resources to cache the weight coefficients between two adjacent layers. The accelerator reads the matrix of weight coefficients data from input buffer, and loops to save into different BRAMs in 32 by the row number of the weight matrix ($n=i\%32$, n refers to the number of BRAM, and i indicates the row number of weight matrix). So, the accelerator can read 32 weight values in parallel. To reduce the impact on performance of the data access time, we design two registers set to read the required data for next iteration computation or be used in current iteration alternately. In our test, the time to cache 32 input values is much less than the time of the computation of 32 values. So, the iteration computing will start without waiting except the first iteration.

C. Adder tree

We use a binary adder tree to implement accumulation. In this way, the time consumption of accumulation can be reduced from $O(n)$ to $O(\log n)$. With the pipeline design, we can get sum or partial sum every clock cycle. When the number of input value cannot be enough (in our case is 32), zeros will be used to be temp input values.

V. EXPERIMENTS AND DATA ANALYSIS

In order to evaluate the performance and cost of the DLAU accelerator, we have implemented the hardware prototype on the Xilinx ZynqZedboard development board, which equips ARM Cortex-A9 processors clocked at 667 MHz and programmable fabrics. For benchmarks, we use the Mnist data set to train the $784 \times M \times N \times 10$ DNNs in MATLAB, and use $M \times N$ layers' weights and nodes value for the input data of DLAU. For comparison, we use Intel Core2 processor clocked at 2.3 GHz as the baseline. In the experiment we use tile size = 32 considering the hardware resources integrated in the Zedboard development board. The DLAU computes 32 hardware neurons with 32 weights every cycle. The clock of DLAU is 200 MHz (one cycle takes 5 ns). Three network sizes— 64×64 , 128×128 , and 256×256 are tested.

Speedup Analysis

We present the speedup of DLAU and some other similar implementations of the deep learning algorithms in Table II. Experimental results demonstrate that the DLAU is able to achieve up to $36.1 \times$ speedup at 256×256 network size.

TABLE II
COMPARISONS BETWEEN SIMILAR APPROACHES

Work	Network	Clock	Speedup	Baseline
Ly&Chow [5]	256×256	100MHz	$32 \times$	2.8GHz P4
Kim et.al [7]	256×256	200MHz	$25 \times$	2.4GHz Core2
DianNao [6]	General	0.98GHz	$117.87 \times$	2GHz SIMD
Zhang et.al [3]	256×256	100MHz	$17.42 \times$	2.2GHz Xeon
DLAU	256×256	200MHz	$36.1 \times$	2.3GHz Core2

In comparison, Ly and Chow [5] and Kim et al. [7] presented the work only on RBM algorithms, while the DLAU is much more scalable and flexible. DianNao [6] reaches up to $117.87 \times$ speedup due to its high working frequency at 0.98 GHz. Moreover, as DianNao is hardwired instead of implemented on an FPGA platform, therefore it cannot efficiently adapt to different neural network sizes.

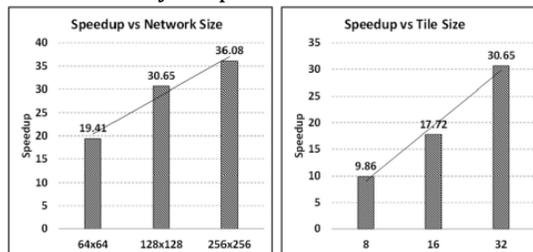


Fig. 5. Speedup at different network sizes and tile sizes.

Fig.5 illustrates the speedup of DLAU at different network sizes- 64×64 , 128×128 , and 256×256 , respectively. Experimental results demonstrate a reasonable ascendant speedup with the growth of neural networks sizes. In particular, the speedup increases from $19.2 \times$ in 64×64 network size to $36.1 \times$ at the 256×256 network size. The right part of Fig. 4 illustrates how the tile size has an impact on the performance

of the DLAU. It can be acknowledged that bigger tile size means more number of neurons to be computed concurrently.

At the network size of 128×128, the speedup is 9.2× when the tile size is 8. When the tile size increases to 32, the speedup reaches 30.5×. Experimental results demonstrate that the DLAU framework is configurable and scalable with different tile sizes. The speedup can be leveraged with hardware cost to achieve satisfying tradeoffs.

Fig. 6 illustrates the floorplan of the FPGA chip. The left corner depicts the ARM processor which is hardwired in the FPGA chip. Other modules, including different components of the DLAU accelerator, the DMA, and memory interconnect, are presented in different colors. Regarding the programming logic devices, TMMU takes most of the areas as it utilizes a significant number of LUTs and FFs.

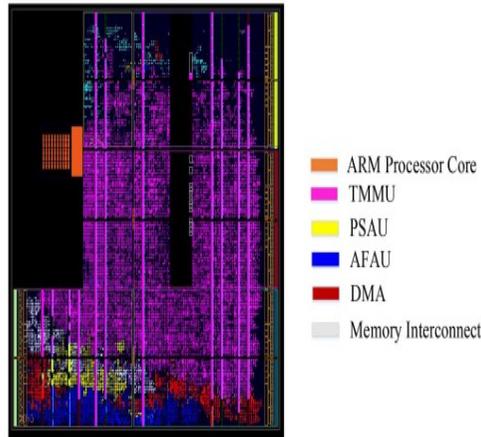
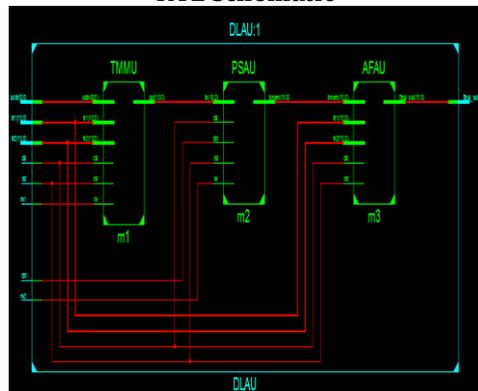


Fig. 6. Floorplan of the FPGA chip.

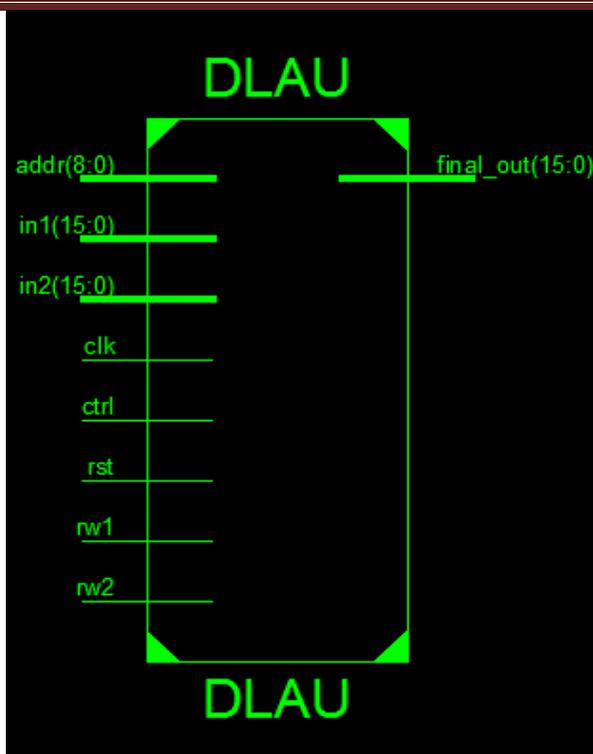
VI.RESULTS
SIMULATION RESULTS



RTL schematic



Technology Schematic:



Design summary

Device Utilization Summary (estimated values)				
Logic Utilization	Used	Available	Utilization	
Number of Slice Registers	658	93120		0%
Number of Slice LUTs	781	46560		1%
Number of fully used LUTFF pairs	578	861		67%
Number of bonded IOBs	62	240		25%
Number of BUFG/BUFGCTRLs	2	32		6%
Number of DSP48E1s	65	288		22%

Timing Report

```

Data Path: m1/m5/Mmult_n0404 to m1/m5/Maddsub_n0408
Cell:in->out      fanout  Gate   Net
                  Delay   Delay Logical Name (Net Name)
-----
DSP48E1:CLK->PCOUT47  1  3.519  0.000  m1/m5/Mmult_n0404 (m1/m5/Mmult_n0404_PCOU...
DSP48E1:PCIN47->PCOUT47  1  1.427  0.000  m1/m5/Maddsub_n0403 (m1/m5/Maddsub_n0403_PCO...
DSP48E1:PCIN47->PCOUT47  1  1.427  0.000  m1/m5/Maddsub_n0406 (m1/m5/Maddsub_n0406_PCO...
DSP48E1:PCIN47      1  1.125  0.000  m1/m5/Maddsub_n0408
-----
Total              7.498ns (7.498ns logic, 0.000ns route)
                    (100.0% logic, 0.0% route)
    
```

VI. EXTENSION WORK

LP-LFSR:

Using Low power test pattern generator, a circuit’s performance can be increased and switching activities can be reduced so that power dissipation will be reduced during test mode. We are going to implement a LP-LFSR techniques so that the power will reduce while speed increases compared to previous works and also switching rate will be reduced in future.

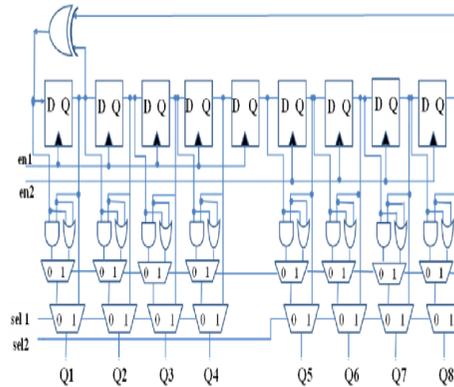
In the proposed LP-LFSR, only one of the flip-flops in the LFSR register is enabled or activated during each shift operation, and as a result the low power LFSR according to the invention generates the same output bit sequence as a conventional LFSR. This lower flip-flop activity causes less power loss and, further, the longer the shift register and the smaller the feedback the more efficient is the low power LFSR.

In order to take number of inputs using less number of clock cycles LP-LFSR is more efficient. Here we use it to take maximum number of inputs from memory in less time which further reduces the delay of the overall system.

Proposed LP-LFSR Architecture:

Fig.7 shows the schematically block diagram of an embodiment of a proposed shift register (LFSR). A low power linear feedback shift register according to the proposed architecture, using memory means such as flip-flops not consuming power when they are disabled. The register does not shift any bits but still

generates the same sequence as a conventional linear feedback shift register. The register comprises enabling means enabling a single current memory means at every shift operation, register stages, each comprising a low power memory means consuming a minimum of power when it is disabled, and feedback means of each stage. Each memory is connected to selection means, selecting at every shift operation the output terminal of a first subsequent memory means following the current memory means being enabled at the current shift operation.



The technique of inserting intermediate vectors is achieved by modifying the conventional LFSR circuit with two additional levels of logic between the conventional flip-flop outputs and the low power outputs as shown in Fig.7. The first level of hierarchy from the top down includes logic circuit design for propagating either the present or the next state of the flip-flops to the second level of hierarchy. The second level of hierarchy is a multiplexer function that provides for selecting between the two states (present or next) to be propagated to the outputs as low power output. Minimal at best consisting of few logic gates.

Only 2 inputs pins, namely test enable and clock are required to activate the generation of the pattern as well as simulation of the design circuit. It is also noteworthy here that the intermediate vectors in addition to aiding in reducing the number of transitions can also empirically assist in detecting faults just as good as the conventional LFSR patterns.

VII. CONCLUSION AND FUTURE WORK

In this paper, we have presented DLAU, which is a scalable and flexible deep learning accelerator based on FPGA. The DLAU includes three pipelined processing units, which can be reused for large scale neural networks. DLAU uses tile techniques to partition the input node data into smaller sets and compute repeatedly by time-sharing the arithmetic logic. Experimental results on Xilinx FPGA prototype show that DLAU can achieve 36.1× speedup with reasonable hardware cost and low power utilization.

The results are promising but there are still some future directions, including optimization of the weight matrix and memory access. Also the trade-off analysis between FPGA and GPU accelerators is another promising direction for large scale neural networks accelerations.

REFERENCES

1. Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
2. J. Hauswald et al., "DjiNn and Tonic: DNN as a service and its implications for future warehouse scale computers," in *Proc. ISCA*, Portland, OR, USA, 2015, pp. 27–40.
3. C. Zhang et al., "Optimizing FPGA-based accelerator design for deep convolutional neural networks," in *Proc. FPGA*, Monterey, CA, USA, 2015, pp. 161–170.
4. P. Thibodeau. Data Centers are the New Polluters. Accessed on Apr. 4, 2016. [Online]. Available: <http://www.computerworld.com/article/2598562/data-center/data-centers-are-the-new-polluters.html>
5. L. Ly and P. Chow, "A high-performance FPGA architecture for restricted Boltzmann machines," in *Proc. FPGA*, Monterey, CA, USA, 2009, pp. 73–82.
6. T. Chen et al., "DianNao: A small-footprint high-throughput accelerator for ubiquitous machine-learning," in *Proc. ASPLOS*, Salt Lake City, UT, USA, 2014, pp. 269–284.
7. S. K. Kim, L. C. McAfee, P. L. McMahon, and K. Olukotun, "A highly scalable restricted Boltzmann machine FPGA implementation," in *Proc. FPL*, Prague, Czech Republic, 2009, pp. 367–372.
8. Q. Yu, C. Wang, X. Ma, X. Li, and X. Zhou, "A deep learning prediction process accelerator based FPGA," in *Proc. CCGRID*, Shenzhen, China, 2015, pp. 1159–1162.
9. J. Qiu et al., "Going deeper with embedded FPGA platform for convolutional neural network," in *Proc. FPGA*, Monterey, CA, USA, 2016, pp. 26–35.