

Multi-ANN embedded system based on a custom 3D-DRAM

Lee B. Baker, Paul Franzon *Fellow, IEEE*

Abstract—Machine Learning in the form of Artificial Neural Networks (ANNs) has gained traction in applications such as image recognition and speech recognition. These applications typically employ a subset of ANNs known as Convolutional Neural Networks (CNNs) which re-use parameters and thus reduce main memory bandwidth. However, there are other types of ANN that do not provide reuse opportunities such as autoencoders and Long Short-term memory. Most research has focused on implementing CNNs but because of their extensive use of SRAM have both ANN size restrictions and performance degradation when used in applications that utilize other types of ANN. This work demonstrates how a customized 3D-DRAM with a very wide databus can be combined with application-specific layers to produce a system meeting the requirements of embedded systems employing multiple instances of disparate ANNs.

Index Terms—machine learning, embedded system, Deep Neural Networks (DNNs), CNN, neural network

I. INTRODUCTION

USEFUL DNNs often require hundreds of thousands of Artificial Neurons (ANes). Within the network, each ANe can have hundreds of feeder (pre-synaptic) ANes. With popular DNNs, there are often tens of layers. So in these ANNs, the memory requirements are significant. The storage is required for the input, the ANe state and most significantly the pre-synaptic ANe connection weights for each of the ANes. This storage requirement often results in gigabytes of memory. When these ANNs are

L. B. Baker, and P. Franzon are with the Department of Electrical and Computer Engineering, North Carolina State University, 2410 Campus Shore Dr., Raleigh NC 27606 Tel/Fax: 919-515-5460/5523 Email: lbbaker@ncsu.edu, paulf@ncsu.edu

Manuscript received Month Day, 2016; revised Month Day, 2016.

required to be solved in fractions of a second, the processing and memory bandwidth becomes prohibitive. In most cases, graphics processing units (GPUs) are used to implement large ANNs. In many ANN architectures, such as Convolutional ANNs (CNN), they are quite effective. However, we should not forget they are not optimized purely for ANN processing and are restricted by available SRAM and they are power hungry. These limitations limit the effectiveness of GPUs.

Much of the ANN Application-Specific Integrated Circuit (ASIC) and Application-Specific Instruction-set Processor (ASIP) research has focused on taking advantage of the performance and ease of use of Static Random Access Memory or SRAM. These implementations can be shown to be effective when there are reuse opportunities such as with CNNs or applications that have batch processing opportunities, such as cloud applications. But given an embedded system requiring multiple disparate ANNs, where reuse opportunities do not exist, these implementations do not provide the required performance. Another technology that has been considered over the last decade is Three-Dimensional Integrated Circuit (3DIC). Implementations such as Neurostream [?], NeuroCube[?] and Tetris [?] employ Hybrid Memory Cube (HMC) 3D-DRAM in a 3DIC system, and all have demonstrated improvements over non-3DIC systems and GPUs. However, for the most part these systems employ the 3D-DRAM as a disjoint memory with the ANN processing elements being an appendage to the dynamic random-access memory (DRAM). In particular, some 3DIC systems that employ 3D-DRAM achieve the required performance by targeting a small process and optimistic scaling numbers where this work achieves the required performance

by tightly integrating the ANN processing with the memory controller while targeting a 28nm process using realistic scaling numbers.

This work considers future embedded systems where multiple ANNs are employed [?] and that each of the ANNs will be of similar size to useful-sized ANNs [?]. In addition, as ANNs fulfill their potential and users start to employ them in various applications, there is no reason to believe these applications will exhibit the same stationary characteristics that make CNNs effective in image processing. Therefore, this work provides support to DNNs that do not provide ANN parameter reuse [?] and suggests that these types of applications will require that all ANN parameters in main memory be accessed in real-time. This work coins the phrase “goldilocks bandwidth” when applied to ANN systems where the system provides the bandwidth required to read all ANN parameters in a real-time sampling rate of 16 ms. Anything exceeding the required bandwidth is considered “too hot” and therefore a waste and anything less than the required bandwidth is “too cold” and therefore inadequate. This work provides support for all types of DNNs including LSTM, autoencoders and restricted boltzmann machines (RBMs) and does not focus on particular ANNs such as CNNs.

Although there is a lot of debate regarding number formats for ANNs, this work also assumes single-precision floating point. This work assumes any useful ANN will be of a similar size to [?], with 772 thousand ANEs and an average fanin to each ANe of 1650. However, the assumption is the convolutional layers employ kernels with non-shared weights. A system employing 10 ANNs for various disparate functions and an average processing time of 16 ms will require an average or “goldilocks” bandwidth of 26 Tbit/s (1) and a capacity of 8.0 GB (2).

By employing a 3DIC system along with customizations to a standard 3D-DRAM, this work demonstrates a system that is able implement multiple useful sized DNNs whilst maintaining the required average (“goldilocks”) memory bandwidth of 26 Tbit/s. This work provides at least a 6X power improvement and a 3X area improvement

Maximum bandwidth

$$\begin{aligned}
 &= \sum_{n=0}^{N_n} \left(\frac{\bar{N}_a \cdot \bar{C}_p \cdot b_w}{\bar{T}_p} \right) \text{bit/s} \\
 &= \sum_{n=0}^9 \left(\frac{772 \times 10^3 \cdot 1.65 \times 10^3 \cdot (32 + 1)}{16 \times 10^{-3}} \right) \\
 &= \sum_{n=0}^9 2.63 \text{ Tbit/s} \\
 &\approx 26 \text{ Tbit/s}
 \end{aligned} \tag{1}$$

where N_n is the number of ANNs

\bar{N}_a is the average number of ANEs

\bar{C}_p is the average number of connections

b_w is the number of bits per parameter

and T_p is the processing time

$$\begin{aligned}
 \text{ANN Memory} &= \sum_{n=0}^{N_n} \left((\bar{N}_p + \bar{N}_a) \cdot b_w \right) \text{Gbit} \\
 &= \sum_{n=0}^9 \left((202 \times 10^6 + 772 \times 10^3) \cdot 32 \right) \\
 &= \sum_{n=0}^9 6.49 \text{ Gbit} \\
 &= 64.9 \text{ Gbit} \equiv 8.1 \text{ GB}
 \end{aligned} \tag{2}$$

where N_n is the number of ANNs

N_p is the number of parameters per ANN

N_a is the number of ANEs per ANN

and b_w is the number of bits per parameter

over similar state of the art systems [?][?][?][?] when those systems are scaled to meet the requirements of these embedded applications. This work removes a reliance on SRAM to achieve high performance thus allowing the proposed design to be utilized in embedded applications when processing multiple disparate ANNs at or near real-time. Although not optimized for specific ANNs, such as CNNs, this work demonstrates the potential for

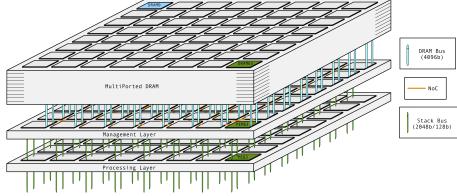


Fig. 1: DRAM, Manager and PE Layers

real-time performance when implementing multiple fully connected DNNs or other similar ANNs such as LSTM.

II. SYSTEM DESCRIPTION

The primary objectives of this work was to consider applications that are unable to take advantage of memory reuse opportunities and therefore not able to achieve high performance using local SRAM, consider embedded devices that will apply many disparate ANNs to perform various system functions, and assume that embedded applications have space and power limitations. This work employs 3DIC technology along with a proposed custom very wide databus 3D-DRAM. The objective was to demonstrate that this very wide databus 3D-DRAM can be used effectively and demonstrate that such a system can stay within the 3DIC footprint of the DRAM. The 3DIC system die stack (figure 1) includes the 3D-DRAM with a system manager below and one or more processing layers below the manager.

3D-DRAM has recently become available in standards such as High Bandwidth Memory (HBM) and Hybrid Memory Cube (HMC) and proprietary devices such as the DiRAM4 available from Tezzaron. These technologies provide high capacity within a small footprint.

In the case of HBM and DiRAM4, the technology can be combined with additional custom layers to provide a system solution.

The question becomes, can a useful system co-exist within the same 3D footprint?

This work targeted a baseline system with:

- single precision floating point for computations
- the Tezzaron DiRAM4 DRAM [?]

The work includes customizing the interface to a 3D-DRAM, researching data structures to describe storage of ANN parameters, designing a memory manager with micro-coded instructions and a processing engine (PE) layer. The targeted 3D-DRAM, the Tezzaron DiRAM4 is a 3D-DRAM employs multiple memory array layers in conjunction with a control and IO layer. The memory is formed from 64 disjoint sub-memories each providing upwards of 1Gigabit with a total capacity of at least 64 gigabit. The system is designed such that a sub-system, known as a sub-system column (SSC) operates on one of the 64 disjoint memories within the 3D-DRAM (see Figure 2).

When the sub-system columns need to share data or neuron activations, the data is passed between SSCs using a mesh connected network-on-chip (NoC).

A control and data block diagram of the 3DIC stack showing the 64 sub-system columns can be seen in Figure 3. A block diagram of the sub-system column can be seen in Figure 4.

An overview of the various blocks and interconnects are given below:

A. 3D-DRAM

The targeted 3D-DRAM, the Tezzaron DiRAM4 is customized to provide a 2048-bit wide bus. A read to the memory using a burst of two cycles provides access to an entire page within a bank. These customizations to support this very wide bus are discussed in IV. The wide bus is connected to the manager using TSVs and the manager directs portions of the wide bus to each lane to the PE.

B. Manager Layer

The Manager block is the main controller in the system. The operations required to process an ANN are formed from individual instructions which are decoded by the Manager. These instructions include descriptors to describe memory read operations, processing engine operations and memory write operations. The manager reads these system

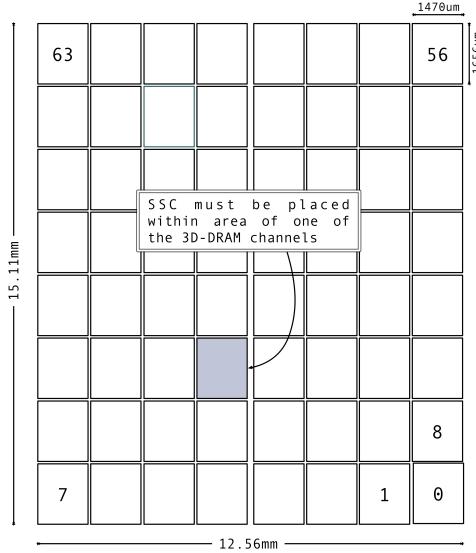


Fig. 2: DiRAM4 DRAM Physical Interface Layout[?][?] showing area for SSC

instructions from an instruction memory, decodes the instruction and configures the various blocks in the system. The configuration includes:

- initiate operand reads from DRAM
- prepare the processing engine (PE) to operate on the operands
- prepare the result processing engine to take the resulting neuron activations from the PE and write those results back to the DRAM
- replicate the resulting neuron activation's to neighbor managers for processing of other ANN layers

C. Processing Layer

The PE is able to operate on data streamed directly from the DRAM via the Manager layer. The PE is configured by the manager to perform operations on the operand data streamed from the manager. In the baseline system, the main operation is to perform multiply-accumulates on 32 execution lanes of two operands. These operands typically are the pre-synaptic neuron activation's and the connection weights. The PE also performs the activation function on the result of the MAC to generate the

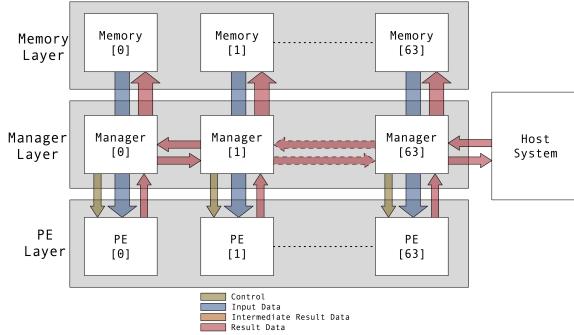


Fig. 3: System Diagram

neuron activation value. These 32 activation values are sent back to the Manager layer.

D. Layer Interconnect

The layers are connected using through-silicon-vias (TSVs) which provide high connection density, high bandwidth and low energy. Figure ?? shows an example of two die connected using TSVs.

E. Inter-Manager Communication

A Network-on-Chip (NoC) allows each management block to communicate with other managers.

During configuration and/or computations, data must be transported between managers. This inter-manager communication is provided by an NoC. When computing an ANN across multiple processors, often neuron activation data must be shared. Each manager contains a NoC module and all the managers are connected in a mesh network. The NoC bandwidth was chosen to ensure the results can be multicast to any destination manager without an adverse affect on the pipelined instructions. The NoC bus width is 64 data bits plus control signals running at the system clock rate. Future testing may require additional NoC bandwidth.

III. SYSTEM OPERATIONS

In the context of this system and ANe state calculation, the basic operations to determine the state of a neuron is to:

- Inform the Manager and PE which operations are to be performed
- Instruct the manager to access the states of the pre-synaptic neurons
- Instruct the manager to access the weights of the connections from the pre-synaptic ANEs
- Provide the pre-synaptic neuron weights and states to the processing engine execution lanes
- Instruct the manager where to store the resulting ANe state back to memory

This work has researched an instruction architecture to describe the above operations. These instructions are decoded by the manager.

The managers primary responsibility is:

- Instruction decode
- Internal Configuration messages
- Operand read
- Result write

The PE has three major blocks:

- Streaming Operation block (StOp)
 - Processes data from the manager on-the-fly without storing in local SRAM
- Single-Instruction Multiple-Data (SIMD)
 - Processes the data from the StOp function
 - * Neuron activation function such as ReLu
 - * Perform non vector operations such as softmax conversion using local SIMD functions, such as e^x and divide
 - Sends the result back to the manager
- DMA/local memory controller
 - Transfer configuration data to PE controller or to store StOp results to a small local SRAM which can be used for access by SIMD or by the StOp function

A detailed block diagram of the SSC can be seen in Figure 4.

IV. SUGGESTED DRAM CUSTOMIZATIONS

This work achieves the increase in bandwidth by proposing that the DRAM expose more of its currently open page.

Without the limitations of having to transfer data beyond the chip stack, this work suggests exposing a larger portion of the page over a 2048-bit wide

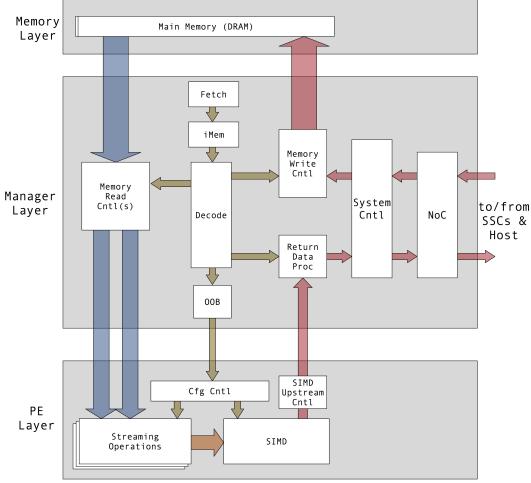


Fig. 4: Sub-System Column (SSC) Block Diagram

bus. By staying within the 3D footprint, this bus can be implemented using fine pitch TSVs.

For every group of ANes processed, the state of the group of ANes is written back to memory. Typically this would require a read/modify/write of a DRAM cacheline. In the case writing back 32 ANe states into a 4096 bit cacheline means the read/modify/write is inefficient. To minimize the inefficiency, a customization to the DRAM is the addition of a write data mask to the DRAM write path eliminating the additional read.

V. RESULTS

The objective of this work was to design a system able to accelerate multiple disparate ANNs in embedded systems. Given that these systems cannot effectively utilize SRAM, the main objective was to demonstrate a system that can operate efficiently using a customized 3D-DRAM with a very wide data-bus. 3DIC technology was exploited, including 3D-DRAM, because a major theme of this work is that 3DIC provides many benefits, including reduced energy use, lower area requirements, and high bandwidth. It was necessary to show that the proposed system can maintain the required data bandwidth while staying within the physical footprint of the 3D-DRAM.

The primary control and datapaths of the system have been simulated in a system verilog environment. Timing closure was at a frequency of 500 MHz and a place and route for the Manager and PE are shown in Figure 5. The parasitics were extracted from these layouts and simulated against a group of operation based on the expected lower and upper limits of pre-synaptic fanin. These testcases were based on layers similar to CONV2 and FC-7 from [?] and represent a pre-synaptic fanin of 225 and 4000 respectively. Additional testcases were employed representing pre-synaptic fanins of 294, 300, 500 and 1000. Both locally connected (CONV) and fully connected (FC) type fanins were tested. The results showing sustained average bandwidth can be seen in Table II.

The estimated overall power along with per block contribution are shown in Table I.

TABLE I: Power Estimates

Technology Node	Clock Frequency	Total Expected Power	Testcase
28nm	500 MHz	64W	CONV-294
28nm	700 MHz	88W	CONV-294

(a) Power Dissipation

Block Name	Percentage Contribution
Manager	66.6 %
PE	28.0 %
DRAM	2.5 %
DRAM TSVs	1.8 %
Stack Bus TSVs	1.2 %

(b) Power Contribution

As bus efficiency is the main metric, Table II shows sustained average bandwidth over the fanin testcases.

VI. CONCLUSIONS

There have been many attempts to accelerate ANNs. Many have shown excellent performance mainly when implementing CNNs. The improvement mostly comes from the ability to hold kernel

TABLE II: Fanin Bandwidth Tests

Test	Average Bandwidth At Frequency	
	500 MHz	700 MHz
CONV2 [?]	~25 Tbit/s	~35 Tbit/s
CONV-294	~26 Tbit/s	~37 Tbit/s
CONV-300	~27 Tbit/s	~37 Tbit/s
CONV-500	~29 Tbit/s	~41 Tbit/s
CONV-1000	~31 Tbit/s	~43 Tbit/s
CONV-2500	~32 Tbit/s	~45 Tbit/s
FC-350	~28 Tbit/s	~39 Tbit/s
FC-500	~29 Tbit/s	~41 Tbit/s
FC-1000	~31 Tbit/s	~43 Tbit/s
FC-7 [?]	~32 Tbit/s	~45 Tbit/s

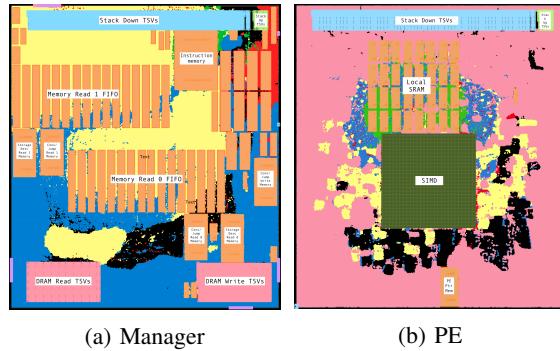


Fig. 5: Manager and PE Die layouts

weights and/or ANe activations in local SRAM. Another method of employing local memory is often due to pooling of batch requests, especially in server applications. This local storage allows the system to take advantage of the low latency and random access benefits of SRAM whilst performing multiple operations on that data. When considering applications where this local storage cannot be used effectively, all these implementations suffer a large degradation in performance.

This work considers embedded applications where a system is processing requests with a disparate set of ANNs. The assumption is that local SRAM is no longer effective and performance is based on DRAM bandwidth. This work considers 3DIC technology and a customized 3D-DRAM is proposed.

The customized 3D-DRAM combined with a design based on custom instructions and operation

descriptors allows the system to achieve high levels of memory bandwidth efficiency.

There is no doubt existing CNN accelerators that take advantage of batch processing achieve a performance that is difficult to better, but applying these systems to this works target application exposes those systems DRAM bandwidth limitations. This work demonstrates a 3DIC system that at the surface provides relatively low floating point operations per second (FLOPS), but considering the target application is memory bound, this work demonstrates a 3X power improvement and 6X area improvement over similar ANN systems [?], [?], [?], [?].

VII. ACKNOWLEDGEMENTS

This work was funded in part by DARPA and AFRL under FA8650-15-1-7518 and DARPA and ONR under N00014-17-1-3013, as part of the CHIPS program.

Lee B. Baker received a B.S. degree in Electrical Engineering from Brighton Polytechnic, UK, an M.S. degree in Electrical Engineering from Villanova University, USA and an M.B.A. degree from North Carolina State University, USA in 1983, 1994 and 2009 respectively. He earned his Ph.D. in Electrical and Computer Engineering from North Carolina State University in 2018. His current research interests include acceleration of artificial neural networks.

Paul Franzon is currently the Cirrus Logic Distinguished Professor of Electrical and Computer Engineering at North Carolina State University. He earned his Ph.D. from the University of Adelaide, Adelaide, Australia in 1988. He has also worked at AT&T Bell Laboratories, DSTO Australia, Australia Telecom and three companies he cofounded, Communica, LightSpin Technologies and Polymer Braille Inc. His current interests center on the technology and design of complex microsystems incorporating VLSI, MEMS, advanced packaging and nano-electronics. He has lead several major efforts and published over 300 papers in these areas. In 1993 he received an NSF Young Investigators Award, in 2001 was selected to join the NCSU Academy of Outstanding Teachers, in 2003, selected as a Distinguished Alumni Professor, and received the Alcoa Research Award in 2005. He served with the Australian Army Reserve for 13 years as an Infantry Solider and Officer. He is a Fellow of the IEEE.