# Algorithms and Hardware for Efficient Processing of Logic-based Neural Networks

Jingkai Hong*, Arash Fayyazi*, Amirhossein Esmaili, Mahdi Nazemi, and Massoud Pedram
Department of Electrical & Computer Engineering, University of Southern California, Los Angeles, CA, USA
{jingkaih,fayyazi,esmailid,nazemi,pedram}@usc.edu

*Abstract*—Recent efforts to improve the performance of neural network (NN) accelerators that meet today's application requirements have given rise to a new trend of logic-based NN inference relying on fixed-function combinational logic (FFCL). This paper presents an innovative optimization methodology for compiling and mapping NNs utilizing FFCL into a logic processor. The presented method maps FFCL blocks to a set of Boolean functions where Boolean operations in each function are mapped to high-performance, low-latency, parallelized processing elements. Graph partitioning and scheduling algorithms are presented to handle FFCL blocks that cannot straightforwardly fit the logic processor. Our experimental evaluations across several datasets and NNs demonstrate the superior performance of our framework in terms of the inference throughput compared to prior art NN accelerators. We achieve 25x higher throughput compared with the XNOR-based accelerator for VGG16 model that can be amplified 5x deploying the graph partitioning and merging algorithms.

## I. INTRODUCTION

Deep neural networks (DNNs) provide state-of-the-art (SoA) performance in various artificial intelligence applications and have surpassed the accuracy of conventional machine learning models in many challenging domains, including computer vision [7], [21] and natural language processing [4], [6]. The emergence of more complex DNN models such as *transformers* [18] and *MLPMixers* [15] is a key reason for the remarkable performance of DNNs in many application domains but these models impose huge compute and memory resource requirements.

To improve the efficiency of DNNs, some prior work formulates the problem of efficient processing of neural networks as a Boolean logic minimization problem where ultimately, logic expressions compute the output of various filters/neurons. NullaNet [10] optimizes a target DNN for a given dataset and maps essential parts of the computation in DNNs to logic blocks, such as look-up tables (LUTs) on FPGAs. In most cases, including those used in this study, the average accuracy drop for binary implementation is less than 4%.

The issue is once the target fixed-function combinational logic (FFCL) for a specific NN model is synthesized, the synthesized fabric can only be used for the inference task of that given NN model, which makes ASIC realization of FFCL impractical. In addition, since neurons designed for SoA NNs include tens to hundreds of inputs, the obtained Boolean logic expression can be huge. Our experiments show that it is impossible to map some generated Boolean logic functions onto a single FPGA because LUTs are used up and cannot accommodate such huge Boolean functions.

A logic processor comprising a Boolean logic unit, in which the Boolean logic unit is responsible for performing Boolean operations of each Boolean function associated with an FFCL block extracted from a binary neural network (BNN), as introduced in [11], is critical to doing inference with different BNNs. Therefore, designing configurable efficient logic processors as logic-based inference engines,

which can be used in a variety of applications with different DNN models, is highly desirable.

The compilation and scheduling of an arbitrary logic graph associated with a Boolean function to be mapped onto a logic processor is a demanding task from the viewpoint of the compiler design. The compiler needs to detect and group the operations of all gates that can be executed simultaneously, considering hardware resource limitations (i.e., the number of Boolean logic units per logic processor). To address these challenges, we design a many-core logic processor and introduce novel techniques for compiling and mapping BNNs that utilize FFCL into this logic processor. Furthermore, we present an original graph partitioning algorithm for handling very large logic graphs. The contributions of the paper are as follows:

- We present the design of a logic processor which can process large Boolean functions.
- We present an innovative optimization methodology for compiling and mapping BNNs utilizing FFCL into this logic processor. The proposed compiler generates customized instructions for static scheduling of all operations of the logic graph during inference.
- We present a method to map FFCL blocks to a set of Boolean functions where Boolean operations in each function are mapped to high-performance, low-latency, and parallelized processing elements.
- We present a graph partitioning algorithm to efficiently deal with very large Boolean logic graphs.
- Our experimental evaluations across several datasets and NNs demonstrate the superior performance of our framework in terms of the inference throughput compared to prior art NN accelerators. We achieve 25x higher throughput compared with the XNOR-based accelerator for VGG16 model [13] that can be amplified 5x deploying the graph partitioning and merging algorithms.

While the scope of this work is focused on the algorithms for enabling the execution of arbitrary-size FFCL blocks on a logic processing fabric, we also present the results of using our algorithms and hardware design on an FPGA to determine their efficacy. The algorithms and hardware can be used for an ASIC design as well.

## II. TERMINOLOGY AND NOTATION

This section provides the terminology and notation used throughout this paper:

- **Fixed-function combinational logic (FFCL) block:** A netlist of a combinational logic circuit written in a hardware description language, such as Verilog.
- **Logic processing element (LPE):** A programmable hardware block which performs (two-input) logic operations such as AND, OR, XOR, etc.

*Jingkai Hong and Arash Fayyazi contributed equally to this work.

- **Logic processing vector (LPV):** A hardware block that contains a fixed number of LPEs. LPVs are linearly ordered relative to each other.
- **Logic processing unit (LPU):** A hardware block that contains a fixed number of LPVs, also called a **logic processor** in this paper.
- **Maximal feasible subgraph (MFG):** A directed acyclic graph (where nodes are Boolean operations and edges are data dependencies) greedily extracted from an FFCL without exceeding the LPU's capacity when mapping to an LPU.
- **Full path balancing (FPB):** Equalizing the logic depth of all propagation paths from circuit inputs (i.e., primary inputs) to circuit outputs (i.e., primary outputs). It guarantees all input-output paths have the same number of gates on them.

## III. Proposed Design Flow

The overall flow of the proposed framework is as depicted in Fig. 1. The input to the flow is a description of an FFCL block in the Verilog language. Please note that the framework can be structured to accept any specification of an FFCL block as the input. Yosys [19] and ABC [3] can be used to generate synthesizable Verilog code from any behavioral specification. NullaNet [10] generates the FFCL block in Verilog format and will be used as the upper stream engine. We first parse the Verilog netlist, synthesize the circuit using standard logic optimization techniques, primarily aimed at reducing the total gate count and depth of the circuit, and map the circuit to a customized cell library. Notice that the Boolean operations supported by the logic gates in the cell library, such as two input AND, OR, and XOR operations, must be supported by the LPEs. Next, the mapped circuit is levelized according to the definition in Section II. The logic synthesis and levelization are the same as the one presented in [12].

Because a gate that is at a specific logic level in a target circuit has no connections to any other gates at the same logic level, operations of all gates at the same logic level can be executed simultaneously. However, these operations may have to be assigned to different compute cycles due to hardware resource limitations, i.e., the fixed number of LPVs per LPU (the *depth issue*) or the fixed number of LPEs per LPV (the *width issue*). Multiple LPUs can be assembled in parallel or series configuration for large graphs to complete the required computations for a given logic graph at the extra area/power cost. Our compiler has the ability to map any logic graph to an arbitrary-size LPU. To handle the width issue, the compiler decomposes the filter/neuron functions into MFGs, each of which is then mapped onto the LPU one after the other (more details in Section V).

## IV. Logic Processor Architecture

The architecture of the logic processor is shown in Fig. 2. The logic processor is a data-driven architecture in the sense that the
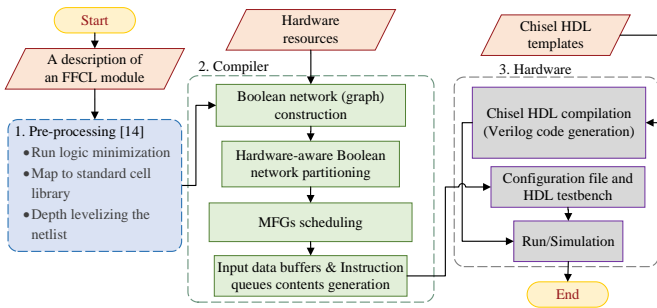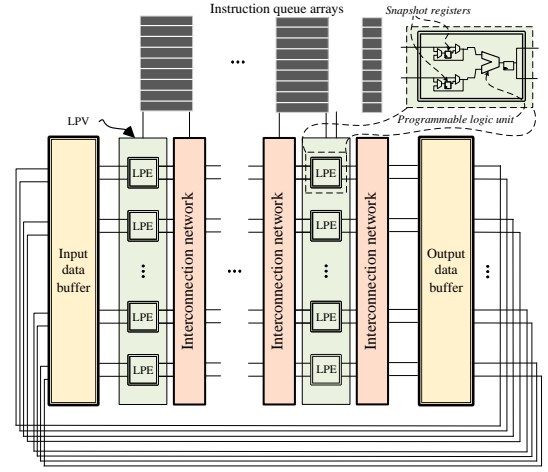


Fig. 2: The LPU architecture.

streaming data is received and processed by multiple stages of LPVs without having to store any intermediate results in some scratchpad memories. More precisely, a logic processor comprises a set of LPVs that are linearly ordered. Each LPV contains $m$ LPEs, each of which receives two inputs and produces one output, resembling a logic gate. Therefore, each LPV receives up to $2m$ input operands and produces a vector of up to $m$ output results. To support logic operation packing and increase hardware efficiency, each operand has a width of $2m$ bits, which translates into $2m$ Boolean variables or $m$ 4-valued logic variables, etc. In the case of processing an FFCL block extracted from a convolutional neural network (CNN) model, the $2m$ bits of data come from different patches of an input feature volume in a CNN or from different images for batch-based inference tasks. To pass data from the $i$th LPV to the $(i + 1)$th LPV, we use a non-blocking multicasting multi-stage switch network.

The number of LPEs per LPV and the number of LPVs per LPU determine the size of the logic graph that can be processed by an LPU. With the parameter values described previously, an LPU can process a logic graph with a maximum width of $m$ and a maximum depth of $n$, where the width refers to the number of logic operations at any logic depth in the graph and the depth refers to the logic depth from any graph inputs to any graph outputs.

Because of the data dependency between two adjacent logic levels within a logic graph, the intuitive way of computing a graph is level by level. However, given the hardware budget is tight and high throughput is desired, it is not an attractive option to use scratchpad memories to store temporary results (i.e., results produced but not yet consumed) between two logic levels. Therefore, we propose distributed *snapshot registers* to store temporary results as shown in Fig. 2. Moreover, we propose a MFG-by-MFG computing paradigm and the scheduling algorithm, which are detailed in Section V. This paradigm further exploits the locality brought about by the *snapshot registers* that are placed in each of the LPEs.

Each LPE contains a *logic unit* where an elementary Boolean operation can be performed, and two *snapshot registers* where each of the LPE inputs can be temporarily stored for a certain data lifecycle determined by the compiler. Additionally, the operations assigned to each LPE are configured with the aid of a instruction set. Notice two types of logic operation can be performed, namely, a multiple-input single-output (MISO) operation including AND, OR, XOR/XNOR and a single-input single-output (SISO) operation including NOT/BUFFER. BUFFER denotes a node that is added to a directed acyclic graph (DAG) to make all paths between any two



Fig. 1: An overview of the proposed framework.

connected nodes have the same topological length (i.e., an equal number of gates exist on all paths between two connected nodes). Buffer insertions are done as a part of the full path balancing step. Full path balancing guarantees no data dependencies exist between two non-adjacent logic levels of gates, simplifying the mapping of the logic graph onto our pipelined architecture.

## V. COMPILER

A compiler tailored to the logic processor is presented. The compiler parses a gate-level Verilog netlist to extract the set of operations that are carried out at each logic level of the circuit netlist, creates a DAG to represent these gate operations and their directional data dependencies, then partitions this DAG into MFGs such that each MFG can fit in the given LPU. Then it schedules the generated MFGs and maps them to different LPVs in different compute cycles. As a result, the LPU processes level by level within each MFG and executes the large logic graph in an MFG-by-MFG manner.

### A. Boolean network partitioning

Consider a subgraph $H = (V', E')$ of a fully path balanced Boolean DAG $G = (V, E)$. Logic levels of nodes in subgraph $H$ are from $L_{bottom}(H)$ to $L_{top}(H)$. Now, given a partitioning solution $P$ of graph $G$ such that $P$ consists of $|P|$ subgraphs, we have:

$$\forall H \in P, \quad 0 \leq L_{bottom}(H) \leq L_{top}(H) \leq L_{max}$$

where 0 and $L_{max}$ denote the levels of primary inputs and primary outputs of $G$, respectively. We denote the set of nodes in a logic level $l$ by $node\_set(l)$. Moreover, given a set of nodes $S$, let $input(S)$ denote the set of distinct nodes that feed into the set of nodes $S$.

The MFGs must satisfy the following conditions. First, the corresponding subgraphs $H$ must satisfy:

$$\forall l \in [L_{bottom}(H) + 1, L_{top}(H)] \implies \tag{1}$$
$$input(node\_set(l)) \in H$$

which implies that inputs of all levels of a subgraph except the bottom-most level must also be contained in the same subgraph. Differently stated, inbound connections from nodes outside the subgraph to a node inside the subgraph are allowed only at the bottom-most level of the subgraph. Second, a subgraph must contain at most $m$ nodes in each level of the subgraph (recall that each LPV in the proposed logic processor has $m$ LPEs):

$$\forall l \in [L_{bottom}(H), L_{top}(H)], |node\_set(l)| \leq m \tag{2}$$

Note that nodes need not belong exclusively to a subgraph, i.e., MFGs can have overlapping node sets. More precisely,

$$\exists l \in H, \exists l' \in H', node\_set(l) \cap node\_set(l') \neq \emptyset \tag{3}$$

And, finally, nodes in the bottom-most level of an MFG must have input node count greater than $m$ with the exception of the subgraphs whose inputs are the PIs of $G$:

$$\forall H \in P \setminus \{H_0 | L_{bottom}(H_0) = 0\}, \\ |input(node\_set(L_{bottom}(H)))| > m \tag{4}$$

where $\setminus$ denotes the set subtraction operation. Notice that the inputs of the bottom-most level of a subgraph $H$ are not included in $H$. This equation is the break condition of while loop in Algorithm 2.

The Boolean network partitioning pseudo-code is shown in Algorithm 1. The algorithm uses a BFS traversal starting from the primary outputs (POs) on the given graph to find all MFGs. The MFG rooted at the POs is obtained by using the $findMFG()$

procedure as explained in Algorithm 2. Algorithm 1 then continues by finding MFGs rooted at input nodes of the just extracted MFG. The traversal continues until we reach the PIs of the Boolean network. The procedure to construct an MFG rooted at some node $V$ keeps adding nodes to the MFG (again in a BFS traversal manner) until it reaches a logic level in its transitive fanin cone that has more than $m$ nodes, which is called the stop level. Note that, as seen in Algorithm 2 and also shown in Fig. 3, the MFG rooted at $V$ does not include the stop level nodes. The MFG cannot include any subset of nodes at its stop level because condition 1 is violated in that case.

Algorithm 1 returns a set of MFGs called allTempMFGs which contains MFGs that have one node at their top level, as shown in Fig. 3. The runtime of a BNN inference task is primarily affected by the total number of MFGs. Therefore, a greedy merging algorithm (see Algorithm 3) is proposed to merge within a set of single-output MFGs that feeds into the same MFG and has the same bottom level, generates one multiple-output MFG (refer to Fig. 3). The checkLevel function checks whether the nodes of two MFGs in each logic level meet the constraint of $|nodes(MFG1, l_i) \cup nodes(MFG2, l_i)| \leq m$. Note that one cannot merge two MFGs that have different bottom levels because of violating the property described in condition 1. Next, we discuss the main challenge we faced in compiler design.

### B. Addressing the width issue

For most of the Boolean networks generated by NullaNet [10], the mapped netlists have sizes exceeding $n$ logic levels and contain a lot more than $m$ nodes per logic level. We refer to the former problem as the depth issue and the latter problem as the width issue, which is addressed by the MFG-by-MFG computing paradigm as follows.

MFGs can be scheduled and processed sequentially with the aid of the *snapshot registers* to store the intermediate results generated by $node\_set(L_{top})$ of each MFG. In a pipelined manner, each MFG requires $L_{top} - L_{bottom} + 1$ number of LPVs for its computation. Precisely, the computational resources allocated to MFG $H$ are $LPV \in [L_{bottom}, L_{top}]$. The compute cycle count is $(L_{top} - L_{bottom} + 1) \times t_c$ clock cycles where $t_c$ is the summation of one cycle for computation within an LPE and $t_{sw}$ cycles for data routing (steering) in the deployed switch network. Note that $t_c$ is fixed for different BNN inference tasks. In this paper, $t_c = 6$ because $t_{sw} = 5$ where a 5-stage non-blocking multicast switch network [20] is used. The *snapshot registers* of LPV associated with the $L_{top} + 1$ of the MFG (e.g., $H_i$) to be routed, store the intermediate results generated by the given MFG. The parent MFG (i.e., $H_j$, the MFG that has inputs from the $H_i$) starts its computation at $L_{top} + 1$ using the data stored in the *snapshot registers*.
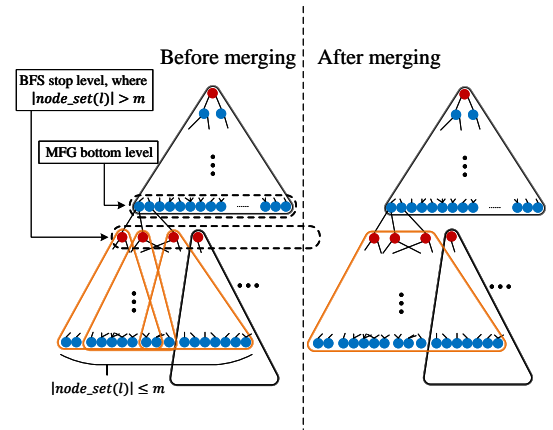


Fig. 3: An example of the merging procedure.

**Algorithm 1** Single-output Boolean network partitioning

**Input**: The PO of a Boolean network, $m$ number of LPEs per LPV
**Output**: A set of MFGs that covers the Boolean network

```
1:  allTempMFGs = []                          // a set of all MFGs
2:  MFG=findMFG(PO,m)                          // call Alg. 2
3:  queue = []
4:  queue.append(MFG)
5:  while queue is not empty do
6:      curMFG = queue.pop(0)                  // pop the first element
7:      allTempMFGs.append(curMFG)
8:      for inputNode in input(curMFG) do
9:          childMFG=findMFG(inputNode,m)      // call Alg. 2
10:         queue.append(childMFG)
11: return allTempMFGs
```

**Algorithm 2** Function *findMFG()*

**Input**: Node $V$, $m$ number of LPEs per LPV
**Output**: A Graph object maintaining a MFG rooted at $V$

```
1:  MFG = Graph(root=V)                        // initialize a Graph with root=V
2:  queue = []
3:  queue.append(V)
4:  visited = emptyset()
5:  visited.add(V)
6:  while queue is not empty do
7:      curNode = queue.pop(0)
8:      MFG.nodeCount[curNode.level]+=1
9:      visited.add(curNode)
10:     if |MFG.nodeCount[curNode.level]| ≥ m then
11:         stopLevel = curNode.level
12:         break
13:     for child in curNode.children do
14:         if child not in visited then
15:             visited.add(child)
16:             queue.append(child)
17: MFG.bottomLevel = stopLevel + 1
18: return MFG
```

**Algorithm 3** MFG merging algorithm

**Input**: A set of MFGs
**Output**: Reduced set of MFGs

```
1:  allMergedMFGs = []
2:  rootMFG = the MFG contained PO(s)
3:  queue = []
4:  queue.append(topMFG)
5:  while queue is not empty do
6:      curMFG = queue.pop(0)                  // pop the first element
7:      allMergedMFGs.append(curMFG)
8:      for all MFG1 and MFG2 in curMFG.children do
9:          merge = False
10:         if MFG1.bottomLevel == MFG2.bottomLevel then
11:             merge = checkLevel(MFG1, MFG2)
12:         if merge then
13:             mergedChildMFG = MFG1 ∪ MFG2
14:             queue.append(mergedChildMFG)
15:             for all grandchild in MFG1.children and MFG2.children do
16:                 update grandchild.parent to mergedChildMFG
17:         else
18:             queue.append(MFG1, MFG2)
19: return allMergedMFGs
```

An illustrative graph partitioning solution is shown in Fig. 4. In this example, a Boolean network is partitioned into 10 MFGs, which are shown with different colors, each is scheduled and processed as shown in the time-space diagram of Fig. 5. The computation periods of different logic levels are shown by $l_i$ after the MFG's name. For instance, $A1$ represents the computation of all nodes in logic level 1 of MFG $A$ and it takes 1 cycle. They are also associated with other two variables, computing cycles (i.e., C $j$) and LPVs (LPV $k$) to represent the computation of $l_i$th logic level of an MFG is performed in LPV $k$ at clock cycle $j$.

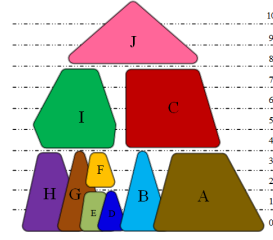In the LPU, a LPV stage and the 5 stages of the subsequent switch
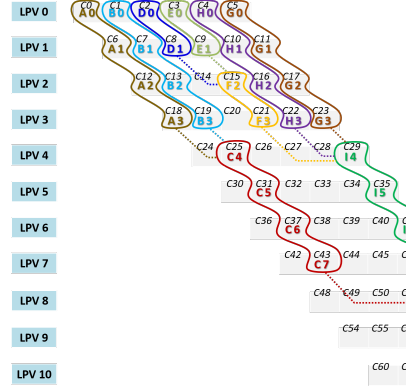


Fig. 4: An example of partitioning where $L_{max} = 10$.



Fig. 5: An example of scheduling. (Note: $C0$ denotes clock cycle 0)

network form a block configured by a 6 instruction queues block, in which each memory takes the read address from its predecessor every cycle. The instruction queues are accessible through a read address shift register. The instructions corresponding to a MFG are written into the same address in all instruction queues associated with all LPVs that are involved in the computations of the MFG. For instance, instructions for computing MFG $J$, highlighted in pink in Fig. 5, are written to memLoc5 of the instruction queues associated with LPVs #8, 9, 10 as shown in Fig. 6.

In Fig. 6, memory locations of instruction queues, which program all the LPVs, are shown. The color of a memory location corresponds to the color of the MFG (cf. Fig. 4). All MFGs with $L_{bottom} = 0$ receive the PI values needed by $node\_set(L_{bottom})$ from the input data buffer. Using a counter, the compiler ensures that the required PI values are properly stored in different locations of the input data buffers such that the desired data is accessed correctly every cycle. This scheme simplifies the address generation compared to a random-access addressing system. Inputs of other MFGs that have $L_{bottom} \neq 0$, are provided by at least two child MFGs that have been computed earlier. Therefore, scheduling the MFGs is equivalent to arranging the instructions and placing them into proper memory locations.

The scheduling algorithm shown in algorithm 4 determines the memory locations for writing the instructions for each MFG. As shown in Fig. 5, each of the MFGs that satisfies $L_{bottom} \neq 0$ has at least two child MFGs supporting all the nodes in its $node\_set(L_{bottom})$. For MFG $H$, we refer to the child MFGs as $H_c = \{H_{c1}, H_{c2}, ..., H_{cn}\}$. Each of the MFGs in $H_c$ will generate a subset of $input(node\_set(L_{bottom}(H)))$; therefore, we have $\bigcup_{i=1}^{n} node\_set(L_{top}(H_{ci})) = input(node\_set(L_{bottom}(H)))$. We refer to the last-scheduled child MFG that generates a subset of the parent's input nodes (without having to store them in the snapshot registers) as the most recent child. For instance, MFG $J$'s most recent child is MFG $I$ (see Fig. 5). Note that two MFGs may use the same memLoc value as long as one of them is the most recent child of the other (e.g., MFG $I$, $J$) because they are performed on different LPVs. So, the required size of the instruction queues is reduced.

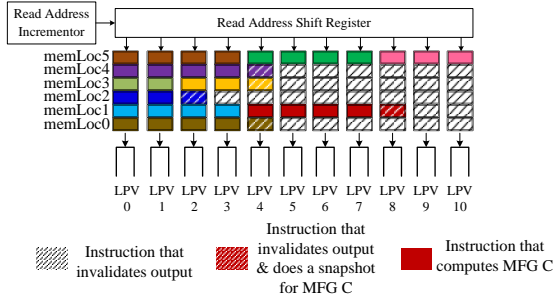Fig. 6: Instruction queue configuration.

**Algorithm 4** MFG scheduling algorithm

**Input**: A set of routed MFGs
**Output**: Memory location to which instructions write

```
 1: memLoc = INT_MAX
 2: topMFG.memLocation() = MemLoc
 3: stack = []
 4: stack.append(topMFG)
 5: while stack not empty do
 6:     curMFG = stack.pop()
 7:     curMFG.memLocation() = memLoc
 8:     if curMFG.parents ≠ None then
 9:         for childMFG in curMFG.children() do
10:             stack.append(childMFG)
11:     else
12:         memLoc = memLoc - 1
13:     for all MFG do
14:         MFG.memLocation() = MFG.memLocation() -
            memLoc
```

### C. Addressing the depth issue

The depth issue is left to the LPU hardware. If an MFG has $L_{top}$ that is greater than the total amount of LPVs of a fixed size LPU, the output data buffer will perform as the snapshot registers of LPV $L_{top} + 1$, which is the LPV that does not physically exist. Instead, LPV 0 resorts to the circulation mechanism to perform the functionality of LPV $L_{top} + 1$. The compiler is responsible for detecting potential depth issue and relocating the instructions accordingly. The compiler notifies the hardware when to feed the intermediate data stored in the output data buffer back to the LPV 0. Such data goes through the pipeline for multiple rounds of computation to complete the computation of all logic levels of the given FFCL.

## VI. EXPERIMENTAL RESULTS

For evaluation purposes, we targeted a high-end Virtex® UltraScale+ FPGA (Xilinx VU9P FPGA , which is available in the cloud as the AWS EC2 F1 instance). We include the FPGA prototyping results since the SoA implementations use the same FPGA. The hardware metrics are reported in table I for LPV count = 16.

Our benchmarked models can be categorized into two groups, i) models for high-accuracy requirement (i.e., large models), and ii) models for high-throughput requirement (i.e., tiny models). In the first group, we study VGG-16, VGG-like model used in ChewBaccaNN [2], LENET5 on the MNIST dataset, and MLPMixers [15] on the CIFAR-10 dataset. We also evaluate our logic processor against extreme-throughput tasks in physics and cybersecurity such as jet substructure classification (JSC) [5] and network intrusion detection (NID) [9]. We used UNSWNB15 dataset to compare our logic processor with other implementations. We employed the same preprocessed training and testing data as that of Murovic et al. [9] which has 593 binary features corresponding to 49 original features and two output classes.

For MLPMixers, the resolution of the input image is 32*32, and the patch size that the experiments are based on is 4*4. So, we have 64 non-overlapping image patches that are mapped to a hidden dimension $C$ which is 128 and 192 for small design (S) and Base design (B), respectively. $D_S$ and $D_C$ are tunable hidden widths in the token-mixing and channel-mixing MLPs, respectively. $D_S$ and $D_C$ are 64 (96) and 512 (768) for S (B) design. There are 8 and 12 mixing layers in S and B designs, respectively.

### A. Effect of the MFG merging procedure

To show the efficacy of the proposed merging procedure described in Section V, we compare the performance of our logic processor with and without incorporating the merging procedure. Results are shown in Fig. 7 and 8. Fig. 7a shows the clock cycle count for computing layers [2:13] of VGG16 with and without the merging procedure. Fig. 7b shows the total number of MFGs obtained from the proposed algorithm with and without the merging procedure. The results confirm the superior performance after applying the merging procedure and also demonstrate the high correlation between computation time and the MFG count. To further assess the effectiveness of the proposed merging algorithm, we apply it to all models used in this study and summarize the results in Fig. 8. As can be seen, the throughput is improved by 5.2x on average while the MFG count can be reduced up to 9.4x.

### B. Comparison Between MAC-based, XNOR-based and NullaNet-based FPGA Implementation of NNs

The VGG16 is a huge network and it has about 138 million parameters. We implement intermediate convolutional layers 2-13 in VGG16 using the proposed framework and fixed-function combinational logic functions. As a baseline for the SoA generic MAC array-based accelerator for the layers realized using conventional MAC calculations, we used the open-source implementation of [14] with some improvements proposed in [12]. We use FINN [16] for our XNOR-based baseline. We improve this implementation by packing operations. We also use the NullaDSP model [12] as another baseline for mapping FFCL generated by NullaNet to the programmable DSP blocks where it can fit any FFCL with any size. We use the best results of each implementation reported in [12] and compare them to our proposed architecture in tables II and III. In the case of JSC and NID, we use the implementation and the associated performance reported in LogicNets [17], Google and CERN's optimized implementation [8], and the implementation presented in [1].

As illustrated in the tables II, our implementation shows superior performance compared to other implementations. The LPU

TABLE I: Resource utilization of design of LPV count = 16.

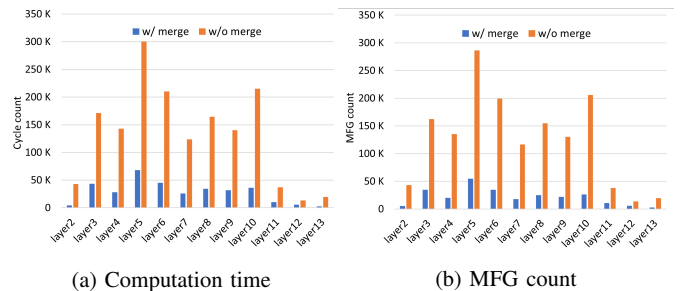| FF(%) | LUT(%) | BRAM(%) | FREQ |
|---|---|---|---|
| 478K(20.2%) | 433K(36.7%) | 12240K (15.8%) | 333MHz |



(a) Computation time    (b) MFG count

Fig. 7: a) Computation time and b) MFG count of VGG16 layers before and after applying merging algorithm.

and XNOR implementation achieves significant saving, especially in large DNN model like VGG16, since we keep all intermediate data on on-chip memories and there is no cost associated with off-chip memories while this is not the case for MAC-based and NullaDSP implementation. We achieved 14.01x(33.43x), 4.86x(3.93x), and 1.95x(4.89x) in performance improvement on VGG16 (LENET-5) inference compared to MAC-based, NullaDSP, XNOR-based implementations, respectively.

LogicNets [17] have higher frames per second (FPS) than our design. However, they cannot use the same hardware for the other models since they realized each model as a customized hard network of logic gates (as in random logic blocks). Whereas, our design offers programmable logic processors that can perform the required logic gate operations of any logic (computation) graphs. The former realization is ideal for building a highly efficient, yet unchangeable, inference engine whereas the latter one is desirable for accelerating the training process and for building inference engines that can be updated after they are deployed in the field. Note that NullaNet Tiny [11] is our upstream for generating FFCL blocks and presents a similar implementation as LogicNets [17] and outperforms the LogicNets in similar settings on the same benchmarks.

*C. Ablation study with LPV count*

To determine the influence of the LPV count on the performance of the presented logic processor, we conducted experiments with different LPV counts. As shown in Fig. 9, the inference time decreases as we increase the number of LPVs. The influence of the LPV count is saturated after a while. To conduct a comparative analysis on NullaDSP [12] against the LPU, we benchmark the effective LPV threshold, which is defined as the minimum number of LPVs that an LPU needs to achieve performance equivalent to NullaDSP. As shown in Fig. 9, we need at least 2 LPVs to achieve such performance for the case of VGG16.

## VII. CONCLUSION

The proposed logic processor offers a novel hardware design approach afforded by the presented compiler resulting in efficient partitioning and mapping of a given neural network model in the format of fixed-function combinational logic. In future work, we plan to intend to explore the heterogeneous architecture where the number
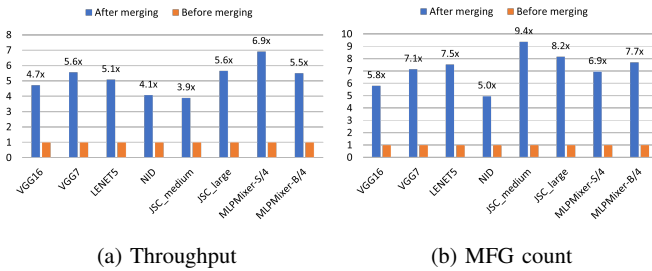


(a) Throughput    (b) MFG count

Fig. 8: a) Throughput and b) MFG count before and after applying merging algorithm.
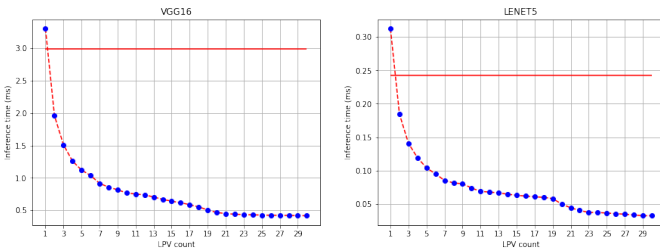


Fig. 9: Inference time of VGG16 and LENET5.

TABLE II: FPS Comparison between the different implementation of models with high accuracy requirements. LPV count in LPU is 16.

|  | MAC | NullaDSP | XNOR | LPU |
|---|---|---|---|---|
| VGG16 | 0.12K | 0.33K | 0.83K | **103.99K** |
| LENET5 | 0.48K | 4.12K | 3.31K | **1035.60K** |
| MLPMixer-S/4 | 4.17K | - | 50.00K | **179.23K** |
| MLPMixer-B/4 | 0.88K | - | 16.67K | **102.01K** |

TABLE III: FPS Comparison between the different implementation of models with high throughput requirements. LPV count in LPU is 16.

|  | LogicNets [17] | Google+CERN [8] | [1] | LPU |
|---|---|---|---|---|
| NID | **95.24M** | - | 49.58M | 8.39M |
| JSC-M | **2995.00M** | - | - | 0.69M |
| JSC-L | **76.92M** | 76.92M | - | 0.21M |

of LPEs per LPVs and their following switch networks will not be the same for all LPVs. Also, it is worth trying multiple LPUs that can be assembled in parallel or series configurations.

## REFERENCES

[1] S. A. Alam et al., "On the RTL implementation of FINN matrix vector compute unit," CoRR, vol. abs/2201.11409, 2022.
[2] R. Andri et al., "Chewbaccann: A flexible 223 TOPS/W BNN accelerator," in IEEE ISCAS, 2021.
[3] R. K. Brayton and A. Mishchenko, "ABC: an academic industrial-strength verification tool," in CAV, vol. 6174. Springer, 2010.
[4] J. Devlin et al., "BERT: pre-training of deep bidirectional transformers for language understanding," in NAACL-HLT. Association for Computational Linguistics, 2019.
[5] J. M. Duarte et al., "Fast inference of deep neural networks in fpgas for particle physics," CoRR, vol. abs/1804.06913, 2018.
[6] S. Hochreiter and J. Schmidhuber, "Long short-term memory," Neural Computation, 1997.
[7] G. Huang et al., "Densely connected convolutional networks," in CVPR, 2017.
[8] C. J. N. C. Jr. et al., "Automatic heterogeneous quantization of deep neural networks for low-latency inference on the edge for particle detectors," Nat. Mach. Intell., 2021.
[9] T. Murovic and A. Trost, "Massively parallel binary neural network inference for detecting ships in FPGA systems on the edge," in IEEE 24th DSD, 2021.
[10] M. Nazemi et al., "Energy-efficient, low-latency realization of neural networks through boolean logic minimization," in ACM 24th ASPDAC, 2019.
[11] M. Nazemi et al., "Nullanet tiny: Ultra-low-latency DNN inference through fixed-function combinational logic," in IEEE 29th FCCM, 2021.
[12] S. N. Shahsavani et al., "Efficient compilation and mapping of fixed function combinational logic onto digital signal processors targeting neural network inference and utilizing high-level synthesis," ACM Trans. Reconfigurable Technol. Syst., 2022.
[13] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in ICLR, 2015.
[14] A. Sohrabizadeh et al., "End-to-end optimization of deep learning applications," in 2020 ACM/SIGDA FPGA, 2020.
[15] I. O. Tolstikhin et al., "MLP-Mixer: An all-mlp architecture for vision," in 34 NeurIPS, 2021.
[16] Y. Umuroglu et al., "FINN: A framework for fast, scalable binarized neural network inference," in 2017 ACM/SIGDA FPGA, 2017.
[17] Y. Umuroglu et al., "Logicnets: Co-designed neural networks and circuits for extreme-throughput applications," in IEEE 30th FPL 2020, 2020.
[18] A. Vaswani et al., "Attention is all you need," CoRR, vol. abs/1706.03762, 2017.
[19] C. Wolf, "Yosys open synthesis suite," 2016. [Online]. Available: http://www.clifford.at/yosys/
[20] Y. Yang and G. M. Masson, "Nonblocking broadcast switching networks," IEEE Trans.s on Computers, vol. 40, no. 9, 1991.
[21] S. Zagoruyko and N. Komodakis, "Wide residual networks," in Proceedings of BMVC 2016. BMVA Press, 2016.