

# Gemmini: An Agile Systolic Array Generator Enabling Systematic Evaluations of Deep-Learning Architectures

Hasan Genc\*, Ameer Haj-Ali\*, Vighnesh Iyer\*, Alon Amid\*, Howard Mao, John Wright, Colin Schmidt, Jerry Zhao, Albert Ou, Max Banister, Yakun Sophia Shao, Borivoje Nikolic, Ion Stoica, Krste Asanovic

University of California, Berkeley

## ABSTRACT

Advances in deep learning and neural networks have resulted in rapid development of hardware accelerators that support them. A large majority of ASIC accelerators, however, target a single hardware design point to accelerate the main computational kernels of deep neural networks such as convolutions or matrix multiplication. On the other hand, the spectrum of use-cases for neural network accelerators, ranging from edge devices to cloud, presents a prime opportunity for agile hardware design and generator methodologies. We present Gemmini<sup>1</sup> - an open source and agile systolic array generator enabling systematic evaluations of deep-learning architectures. Gemmini generates a custom ASIC accelerator for matrix multiplication based on a systolic array architecture, complete with additional functions for neural network inference. Gemmini runs with the RISC-V ISA, and is integrated with the Rocket Chip System-on-Chip generator ecosystem, including Rocket in-order cores and BOOM out-of-order cores. Through an elaborate design space exploration case study, this work demonstrates the selection processes of various parameters for the use-case of inference on edge devices. Selected design points achieve two to three orders of magnitude speedup in deep neural network inference compared to the baseline execution on a host processor. Gemmini-generated accelerators were used in the fabrication of test systems-on-chip in TSMC 16nm and Intel 22FFL process technologies.

## 1. INTRODUCTION

Deep neural networks [1] (DNNs) have gained major interest in recent years due to their extraordinary and robust ability to make predictions on large amounts of data. These prediction abilities have been applied to computer vision [2], machine translation [3], gaming [4], robotics [4, 5], and many other fields. Hardware accelerators are a natural solution to the large computational requirements imposed by DNNs.

A large portion of DNN accelerators produced by major vendors such as Google [6], Samsung [7] and Tesla [8] have used systolic array architectures for matrix multipli-

cation and convolution operations. Systolic arrays were originally proposed in the 1980s [9, 10], but have recently regained interest from their effectiveness in accelerating general matrix multiplications (GEMM) and convolutions in modern machine-learning (ML) workloads.

Accelerators can be used in various stages of the machine learning process: whether in training or inference, on edge devices or on the cloud. Each of these use cases applies different constraints on the accelerator, including latency, power, throughput, energy, area, programmability, and system integration. Nevertheless, the intrinsic computational kernels used in these scenarios remain the same. Critically, the differences between edge inference and cloud training accelerators can be cast as different accelerator parameters rather than changes to the basic computational kernels.

For these reasons, hardware generators [11, 12] are an attractive approach to building DNN accelerators. Although computational kernels may stay the same across workloads, characteristics such as layer dimensions or model size impact how workloads are optimally scheduled and mapped to any particular hardware accelerator [13]. Thus, full-stack generators must target software frontends as well as hardware backends, so that workloads and accelerators can be tuned together.

Systolic array hardware generators should target pertinent architectural parameters such as dataflow, pipeline depth, banking strategy, precision, and on-chip memory capacity. Such generators also need to consider parameters for system-level integration such as bus widths, off-chip memory bandwidth, and host CPU architecture. Accurately evaluating the generated system requires a high-fidelity simulator which can faithfully model system-level interactions with operating systems, DRAM controllers, networks, etc.

Many of these architectural parameters impact the physical realizability, as well as the power, area, and maximum clock frequency of the generated hardware. Therefore, any generator needs to be evaluated not only on its architectural or RTL characteristics, but also on the full physical design flow which it enables.

In this paper, we address these needs and present Gemmini, an agile systolic array generator, which is integrated with the Rocket Chip system-on-chip generator [14] and the BOOM out-of-order processor [15].

\*Equal contribution.

<sup>1</sup><https://github.com/ucb-bar/gemmini.git>

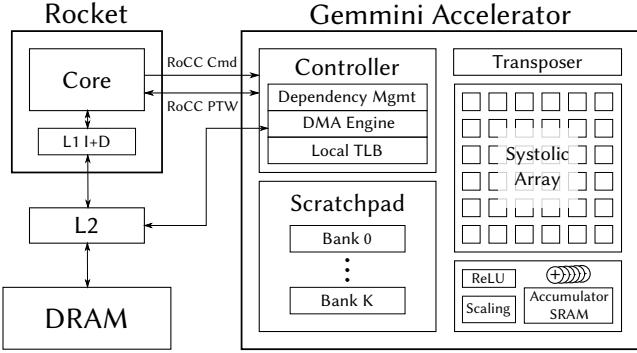


Figure 1: A system overview of the Gemmini based systolic array generator.

Gemmini is composed of a hardware/software generator which produces both RTL and optimized C libraries for common neural network (NN) workloads. We utilize Firesim [16], a cycle-exact FPGA-accelerated simulation platform, to extract accurate performance figures from a design-space exploration (DSE) over architectural and system-level parameters. Additionally, we evaluate Gemmini across the full physical design flow to produce tapeout-ready designs with varying timing constraints and floorplans.

Our DSE revealed that evaluation of any one neural network layer in isolation is not representative of performance on an entire network, partly because the performance and energy-efficiency of different layers can vary widely based on their dimensions and tiling factors. Furthermore, we show that CPU performance and system-level constraints such as bus protocols can severely limit the maximum performance an accelerator can provide. We also demonstrate that Gemmini can produce tapeout-ready designs which meet timing and power constraints with a variety of different floorplans and configurations. Gemmini designs have even been taped-out in TSMC 16nm and Intel 22FFL process technologies.

## 2. GEMMINI GENERATOR

Gemmini is an open-source modular and flexible generator of systolic array accelerators, supporting multiple dataflows, targeting ASIC and FPGA implementations. Gemmini is open source, written in the Chisel hardware description language [17], enabling parameterization and configurability through high-level meta-programming and functional programming abstractions. Gemmini produces instances of systolic architectures that can be integrated with the Rocket Chip SoC generator. Its parameterization and system-level integration enable efficient hardware and software co-design, and help perform agile design space exploration. This section describes the architecture of a systolic array generated by Gemmini (Section 2.1), the major generator parameters (Section 2.2), and the accelerator programming model (Section 2.3).

### 2.1 Architecture

A system-level view of the Gemmini generated accelerator is illustrated in Figure 1. The core unit is a 2-D systolic array that performs matrix multiplications, represented by the equation:

$$C = A * B + D$$

where  $A$  and  $B$  are the multiplied matrices,  $C$  is the result and  $D$  is a bias matrix. The array is fed by a banked scratchpad memory made of SRAMs, with access to main system handled by a direct memory access (DMA) engine in the controller. There are dedicated components for non-linear activation functions, such as ReLU and ReLU6, as well as components necessary for retaining network accuracy after quantization [18], such as rounding and saturating bitshifts. The accelerator also includes an accumulator with a wider bitwidth than the systolic array itself to accumulate partial results.

The  $\mu$ Arch of the systolic array is illustrated in Figure 2. The basic element of the systolic array is a fully combinational processing element (PE), which performs MACs, and optionally rounding bitshifts. The PEs can support a variety of dataflows, which may either be fixed at design time or configurable at runtime. The PEs can support different bitwidths for their inputs, outputs, and internal buffer (Section 2.2), as determined at elaboration time. To enable full utilization of the MAC units, each PE is double-buffered such that weights/biases can be loaded for a future computation while the current compute cycle is running. PEs are arranged in a combinational grid to form a *tile*, and tiles are arranged in a pipelined grid to form the systolic array itself.

To perform a GEMM operation,  $A$ ,  $B$ , and  $D$  matrices must be explicitly moved into the scratchpad from main memory ( $D$  may also be moved directly into the accumulator). The systolic array is then configured with the desired dataflow and activation functions. Afterwards, the  $A$ ,  $B$ , and  $D$  matrices are fed directly into the systolic array, which writes the result,  $C$ , either back into the scratchpad or into the accumulator. Finally, the result may be written back into main memory.

For workloads that are sensitive to precision and rounding, the result of a matrix multiplication must often be of a higher bitwidth than the input matrices. To support this pattern, the Gemmini architecture also includes a higher-bitwidth accumulator external to the systolic array, which is implemented as a dual-port SRAM with adders at its inputs.

The template architecture also includes peripheral circuitry, which performs activation functions and scales high-bitwidth values down to lower-bitwidth values if necessary. For example, Gemmini supports rounding bitshifts, which can be applied within PEs (for the output-stationary dataflow) or at the output of the accumulator (for the weight-stationary dataflow). In a quantized neural network, output activations are usually accumulated to higher precision, *e.g.*, 32 bits. However, before being fed into other layers, these activations must be scaled back down to a lower precision, such as 8 bits. Gemmini saturates and rounds such scaling operations to the

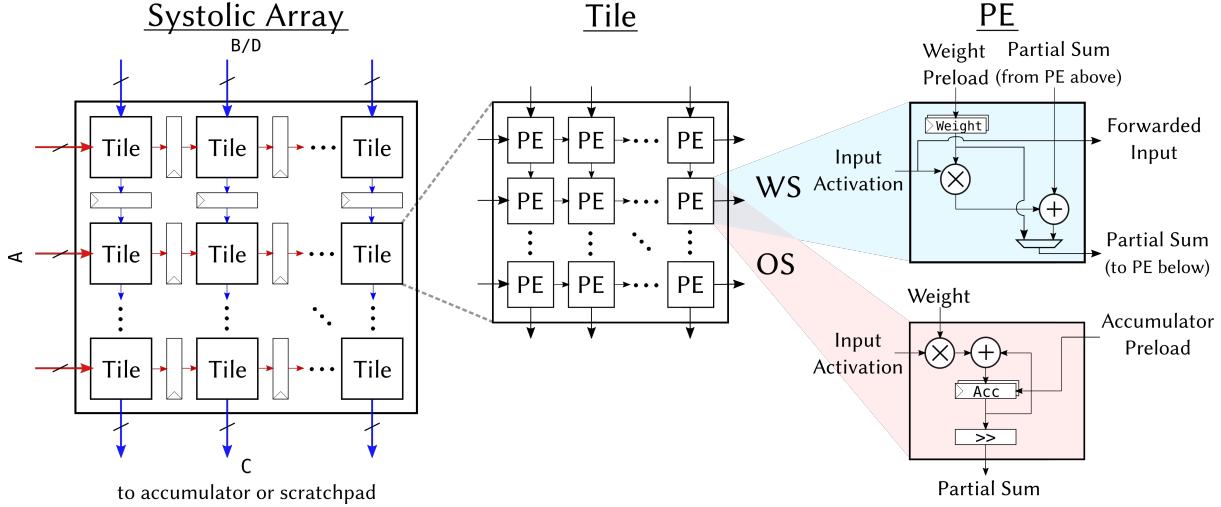


Figure 2:  $\mu$ Arch of the Gemmini systolic array. PEs are arranged in a combinational grid to form a tile, and tiles are arranged in a pipelined grid to form the systolic array itself

nearest bit in order to maximize accuracy [18].

Some of this peripheral circuitry also preprocesses the data. For example, our architecture includes a dedicated transposer, which is itself implemented as a smaller systolic array. For the output-stationary dataflow, a PE must consume the *rows* of  $A$  while consuming the *columns* of  $B$ . However, typically all matrices are stored in main memory as row-major. The transposer allows both  $A$  and  $B$  to be stored as row-major and makes matrix transformations transparent to the programmer.

The accelerator is integrated with the Rocket Chip System-on-Chip generator, which can be configured to use either the Rocket [14] in-order core or the BOOM [15] out-of-order core. The accelerator communicates with the host processor through the Rocket Co-Processor (RoCC) interface, which enables the host RISC-V core to send the accelerator a stream of custom instructions. The RoCC interface enables the accelerator to be integrated into the processor’s cache-coherent TileLink [19] memory system, and provides execution ordering semantics with respect to the host processor.

## 2.2 Parameters

Although all the accelerator instances produced by Gemmini have the same general architecture, a designer can explore different trade-offs in performance, energy and area, based on a range of tunable generator parameters. Choosing an appropriate design point for a specific application is extremely important in the case of a general kernel such as matrix multiplication. As such, an energy-conscious accelerator for a mobile device with limited parallelism would likely choose smaller array sizes and a single dataflow (at the cost of performance and flexibility), while larger cloud-based accelerators with batch-level parallelism can choose larger array sizes and multiple dataflows (for optimal performance). Some of the current parameters enabled by Gemmini are described below.

**Dataflow:** a dataflow describes the data movement

into and out of the systolic array and the communication patterns between PEs. In the classic three-level nested for-loop for matrix multiplications, the dataflow determines which loops are unrolled spatially and which are unrolled temporally. Currently our generator supports both the output-stationary and the weight-stationary dataflows. The dataflow can either be fixed at elaboration time (improving energy efficiency and physical design), or configured at runtime (improving flexibility and possible performance). Previous work has demonstrated that runtime configurable dataflows can improve DNN inference performance and energy efficiency [13].

**Dimensions:** systolic arrays can be generated with any number of PEs the user chooses. As arrays get larger, more operations will be executed per cycle, and data reuse will improve. However, large arrays increase the latency of small matrix multiplications, as operands must traverse the entire length and height of an array before the result can be read out. Large arrays also suffer from low utilization when operating on small matrices, wasting energy and area. Furthermore, large arrays can have a significant impact on physical design and cycle time, since the scratchpad memories need to be placed appropriately to reduce wire-delay between the memory and the array edges.

**Bitwidth:** the generator can be configured at elaboration time to operate on arbitrary bitwidth matrices. The final accumulated result of a matrix multiplication can also have a different bitwidth than the input matrices. Previous work has demonstrated that DNN compression and quantization enable significant energy savings [20, 21] at the potential cost of accuracy. Our bitwidth parameterization enables a designer to explore the accuracy-efficiency trade-off and choose an appropriate design point for the respective application.

**Pipeline Depth:** traditionally, systolic arrays place registers between each PE. However, our generator allows the density of these registers to be reduced, even to the

point of the array being made of fully combinational logic. Fewer pipeline registers reduce the area requirement for our accelerator, but may reduce the maximum achievable clock frequency. The optimal pipeline depth is impacted by physical design and the choice of fabrication process technology.

**Memory Capacity:** both the scratchpad and accumulator memories (implemented using SRAMs) can be configured to have arbitrary capacities. Previous work has found that data movement and coherency management between main memory and accelerators' private memory can consume up to 40% of an accelerated workload's total runtime [22]. Since data-transfer between main memory and the private scratchpad/accumulator memory is expensive, it is beneficial to have large scratchpads to allow for maximal data re-use. However, over-provisioned private memory can lead to energy and area inefficiency. Therefore memory capacity should be balanced with the system bus and DMA bandwidths which drive the memory as well as the data re-use potential of the accelerated workload.

**Memory Banks:** the private memory scratchpad is divided into banks in order to maximize read/write throughput. A larger number of banks allows for higher throughput, but results in additional wiring and physical design constraints.

**System Parameters:** since Gemmini is integrated with the Rocket Chip SoC ecosystem, it can use SoC-level parameters which have been shown to have an impact on accelerator performance [22]. One such parameter is the host processor, which can be an in-order Rocket core or an out-of-order BOOM core. Another example is the SoC system-bus width, which impacts the bandwidth with which the accelerator can communicate and move data between main memory and the private scratchpads.

**Datatype Parameters:** through the use of Chisel hardware description language and Scala typeclass features, our generator is type-generic over the concrete datatype being processed by the systolic array. Gemmini can create accelerator instances which operate on signed integers, unsigned integers, floating point values, or any user-defined datatype, such as a posit [23] or dynamic fixed-point number, through the implementation of the relevant Scala typeclass. This level of parameterization can enable the generator to produce instances specialized for low-precision DNN integer inference operations, as well as for high-precision floating point DNN training and scientific computing.

### 2.3 Programming Model

Gemmini is programmed via a stream of custom RISC-V instructions transmitted directly from a host processor to our accelerator. Gemmini connects directly to the datapath of a RISC-V core, through the Rocket Custom Coprocessor Interface [14]. The accelerator has its own instruction queues, allowing it to run in parallel with the host processor.

Data and memory management between the accelerator and the host processor is explicit, *i.e.*, data must

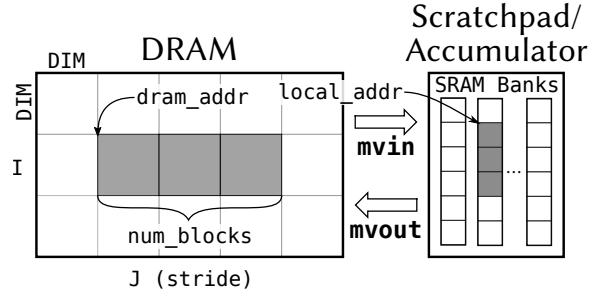


Figure 3: Moving matrices in and out of private scratchpad or accumulator memory with the `mvin` and `mvout` instructions.

be explicitly moved between the processor's main address space and the accelerator's private address space using a sequence of movement instructions. The ISA defines two data movement instructions `mvin` and `mvout` shown in Figure 3. These instructions use Gemmini's DMA unit to move multiple systolic-dimension (DIM) matrices between main memory and the accelerator's private memory space consisting of the scratchpad and accumulator's SRAM banks.

Once matrices have been brought in from main memory, the Gemmini ISA provides a compute instruction that can be configured with a dataflow, a scaling factor, and an activation function. The `compute` instruction takes the local addresses of the  $A$ ,  $B$ ,  $C$ , and  $D$  matrices which can be stored in any scratchpad or accumulator bank.

The output stationary (OS) variant of `compute` (illustrated in Figure 4) executes by loading the  $D$  matrix into the PEs' internal accumulators, pushing  $A$  and  $B$  through the systolic array, and leaves the result  $C$  resident in each PEs accumulator. Providing addresses for the  $D$  and  $C$  matrices are optional in the OS case. This is useful, for example, when a programmer wants to repeatedly accumulate submatrix multiplications on top of each other without reading the results out of the systolic array until the final result has been calculated.

The weight stationary (WS) variant of `compute` (illustrated in Figure 5) takes local addresses for  $A$ ,  $B$ , and  $C$ . First,  $B$  is preloaded into the PEs' weight buffer, then  $A$  is pushed through the systolic array, and the result  $C$  is written to the accumulator. A bias matrix  $D$  can be used in the WS dataflow by first executing a `mvin` into the accumulator. Specifying  $B$  is optional, so the programmer can reuse the already loaded weights in the systolic array.

The Gemmini architecture uses a decoupled-access-execute [24] architecture, where all instructions are issued to one of three independent, parallel command queues: the LOAD queue (`mvin`), the STORE queue (`mvout`), and the EXECUTE queue (`compute`). Any data hazards *within* a command queue are handled transparently by hardware. However, dependencies *between* queues must be encoded into the instructions themselves by the compiler or programmer. Each instruction has four reserved bits which specify whether the instruc-

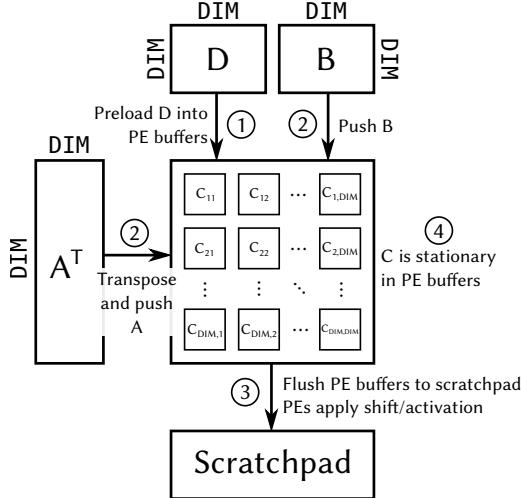


Figure 4: Execution of an output-stationary compute instruction.

tion depends upon an instruction in another queue, or whether an instruction in another queue will depend upon it. This scheme is inexpensive to implement in hardware, but it increases software complexity. Similar software-based dependency management schemes have been implemented in other NN accelerator works [25].

To make it easier for programmers to use Gemmini accelerators, we provide a software library that implements hand-tuned, tiled GEMM functions (supporting both dataflows): matrix multiplication of any size, multi-level perceptrons (MLP), convolutional neural networks (CNN), non-linear activations, and quantization. Tiling is performed along the parameterized size of the systolic array and the accelerator scratchpad. The tiling parameters are generated by the Chisel generator, and included as a header file in the software libraries. This approach facilitates rapid software-hardware co-design.

### 3. DESIGN SPACE EXPLORATION

A major advantage of a generator-based methodology is the ability to perform elaborate design space exploration across a multi-dimensional design space. In this section we explore the performance of multiple design points on different DNN workloads.

#### 3.1 Evaluation Method

We chose to run the DNN applications under Linux to evaluate their performance in a full-system context. Performing this type of RTL evaluation using a logic simulator would take multiple compute-years. Therefore, for full-system performance evaluation we used FireSim, an FPGA-accelerated cycle-exact simulation platform [16]. Unlike FPGA prototyping, FireSim is designed to simulate ASIC RTL designs with timing-accurate system components. FireSim facilitates full-system simulation by enabling integration of the simulated SoC with accurate peripheral and system-level interface models such as DDR3 memory and a last-level-cache (LLC) [26]. By

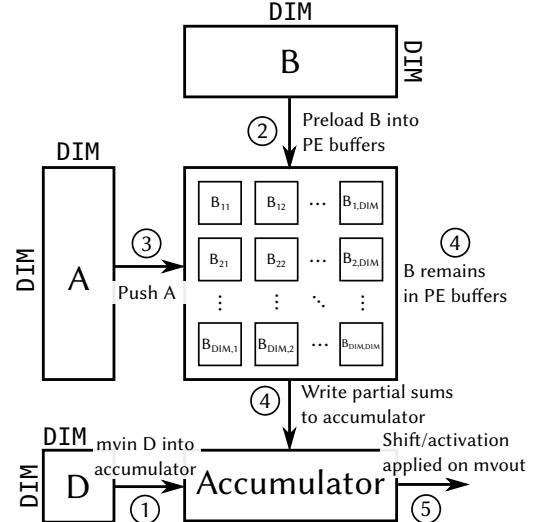


Figure 5: Execution of a weight-stationary compute instruction.

using FireSim’s timing-accurate models, we faithfully simulate our target ASIC designs.

Power and area are evaluated using a Cadence VLSI flow with TSMC 16 nm FinFET technology libraries. Logic synthesis was performed using Genus, physical design was performed using Innovus, and power estimation was performed using Voltus. The accelerators were synthesized to meet frequencies of 1 GHz and 500 MHz.

For performance evaluation, the memory system includes a 256 KiB L2 Cache and a 4 MiB last level cache (LLC). The simulated backing memory preserves the same latency for the 500 MHz and 1 GHz design points, while proportionally scaling the memory bandwidth for the 1 GHz design. At 500 MHz, we used the DDR3 1066 8-8-8 model and at 1 GHz we used the DDR3 2133 14-14-14 model.

Each design point for the DSE was selected by varying a single design parameter relative to a *baseline* which matches design point ① in Table 1. This method attempts to identify and isolate the impact of the different parameters on area, performance, and power consumption. The baseline design point was selected based on common parameters published in the literature.

#### 3.2 Area and Power

The area and power consumption of our designs, normalized and illustrated in Figure 6, were heavily correlated. According to synthesis our baseline 500 MHz design took up 0.467 mm<sup>2</sup> and consumed 611 mW, when including the area and power of the RISC-V core connected to our systolic array. At 1 GHz, synthesis reported that the baseline design took up 0.541 mm<sup>2</sup> and consumed 1.4 W. However, we found that synthesis results were generally pessimistic. Designs which we place-and-routed sometimes consumed less than half of what synthesis predicted they would consume, as seen in Section 4. The trends between different design points, however, remained the same.

Table 1: Design Points Under Evaluation.

No.	Dataflow	Bitwidth	Dimensions	Pipeline Depth	Memory	Banks	Bus width	Host CPU
①	OS	8 bit input 32 bit result	$16 \times 16$	fully pipelined	64 KiB	5	128 bits	rocket
②	<b>WS</b>	8 bit input 32 bit result	$16 \times 16$	fully pipelined	64 KiB	5	128 bits	rocket
③	<b>OS + WS</b>	8 bit input 32 bit result	$16 \times 16$	fully pipelined	64 KiB	5	128 bits	rocket
④	OS	<b>32 bit input 32 bit result</b>	$16 \times 16$	fully pipelined	64 KiB	5	128 bits	rocket
⑤	OS	8 bit input 32 bit result	<b><math>32 \times 32</math></b>	fully pipelined	64 KiB	5	128 bits	rocket
⑥	OS	8 bit input 32 bit result	$16 \times 16$	<b>fully combin.</b>	64 KiB	5	128 bits	rocket
⑦	OS	8 bit input 32 bit result	$16 \times 16$	fully pipelined	<b>256 KiB</b>	5	128 bits	rocket
⑧	OS	8 bit input 32 bit result	$16 \times 16$	fully pipelined	64 KiB	<b>33</b>	128 bits	rocket
⑨	OS	8 bit input 32 bit result	$16 \times 16$	fully pipelined	64 KiB	5	<b>64 bits</b>	rocket
⑩	OS	8 bit input 32 bit result	$16 \times 16$	fully pipelined	64 KiB	5	128 bits	<b>BOOM</b>

Each design point varies a single parameter compared to the baseline (①). The “Banks” column describes the number of scratchpad banks with the addition of an additional bank for the accumulator (which is in its own memory address space).

The weight-stationary dataflow (②) consumed less power than the output-stationary baseline, as it did not require 32-bit accumulators in the PEs of the systolic mesh. Configurations which increased the size of the systolic mesh, on the other hand, such as by scaling up its dimensions or bitwidth, increased power consumption by up to  $3.4\times$  and area by up to  $2.3\times$ . Design (⑩), which replaced the default in-order Rocket processor with a four-wide out-of-order BOOM processor also significantly increased both area and power consumption, whereas in the other design points, the CPU had only a minor impact upon the overall power and area.

### 3.3 Performance

We evaluate the selected design points by running DNNs such as MobileNet, ResNet50, and Resnet152, as well as an additional collection of MLPs, which we refer to in Figure 7 as MLP 1 [27], MLP 2 [28], MLP 3 [29], MLP 4 [30]. The evaluated DNNs represent a wide range of modern state-of-the-art neural network architectures. They include MLPs (which make up more than 61% of Google’s inference workloads [6]), autoencoders, non-linear activations, convolutions, quantization, and depthwise convolutions.

We observed that many of the design points that were expected to boost performance did not have a large impact due to system-level and  $\mu$ Arch effects, while also noting significant variability in the performance boosts achieved by different workloads.

As seen in Figure 7a, for DNN workloads, using a beefier processor in (⑩) boosted performance substantially, while increasing scratchpad memory had little impact, contrary to typical intuition. Since the DNN workloads used the CPU core to perform tasks that map poorly to GEMMs, the CPU often became the bottleneck that limited the maximum speedup achievable. For example, our DNNs performed `im2col` reshaping [31,32] to convert 2D convolutions to GEMMs. With MobileNet in particular, accelerated computation time was dominated by depthwise convolutions on the CPU. Some layers of the evaluated DNNs include  $1 \times 1$  convolutional kernels that could be mapped directly to matrix multiplication without requiring any reshaping. Resnet-152 included the highest portion of such kernels, and thus it performed better in general in all the design points.

While the larger scratchpad (⑦) added more data locality, improving performance by a marginal  $1.18\times$  on our DNNs, its benefit was limited by the CPU bottleneck. On the other hand, increasing bitwidths to 32 bits (④) reduced performance significantly in all cases, as it caused memory requirement to increase, limiting re-use and locality within the scratchpad.

For MLP workloads, the CPU is only used for bookkeeping, so increasing the memory and compute capacity of the accelerator had a larger impact on performance as seen in Figure 7b. Increasing the host CPU’s performance did help, but not as substantially as increasing the dimensions of the systolic array or boosting its scratch-

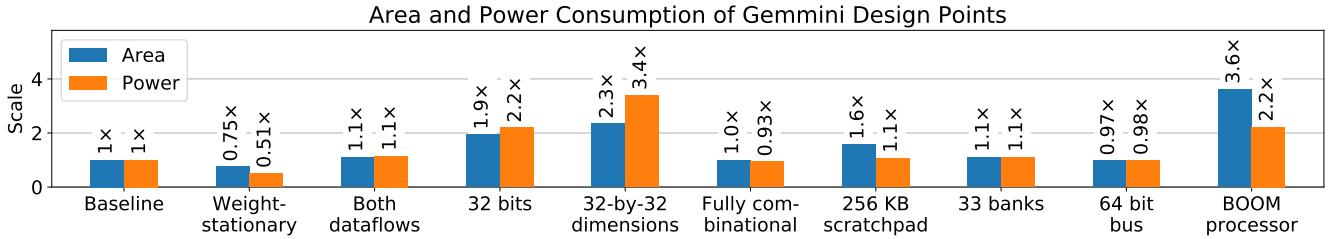
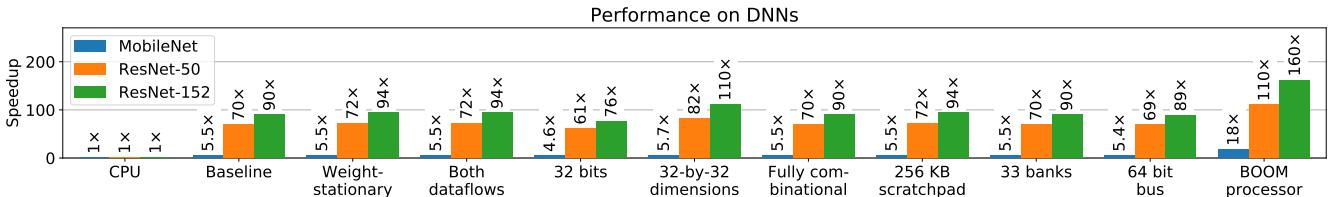
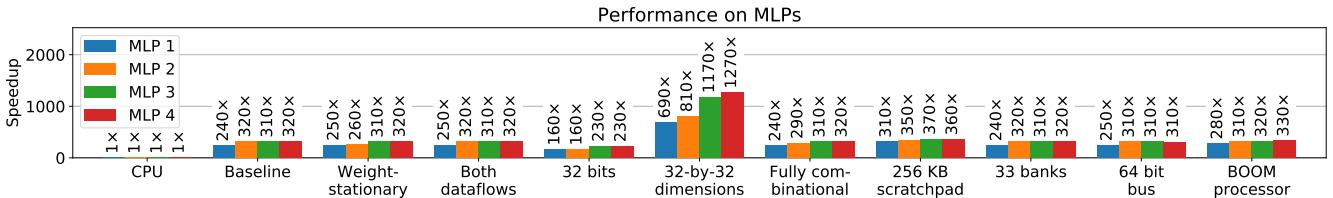


Figure 6: The area and power consumption of synthesized Gemmini designs, normalized to the area and power consumed by the baseline design.



(a) The performance of Gemmini designs on various deep neural networks, normalized to the performance of a cache-blocking algorithm on a CPU.



(b) The performance of Gemmini designs on various multi-layer perceptrons, normalized to the performance of a cache-blocking algorithm on a CPU.

Figure 7: The performance of different Gemmini design points when performing matrix multiplications neural network inference.

pad size.

One would expect that Gemmini is memory-bandwidth limited, and thus cutting the memory bus width would degrade performance. However, we observe no significant performance hit in design ④ owing to a system-level limitation on the number of memory requests in-flight. This limitation turns a bandwidth constraint into a memory latency constraint. Since the round-trip latency of a memory request and the number of maximum requests in flight are independent of the bus width, decreasing it does not impact the effective bandwidth. This reveals the critical importance of system-level evaluation, since using an ideal memory model at Gemmini’s memory port would not reveal system-level bottlenecks.

We observe up to 4 $\times$  performance improvement on MLP inference when increasing the size of the systolic array to 32  $\times$  32 (⑤). The Gemmini  $\mu$ Arch requests multiple systolic-dimension matrix rows at a time when executing the `mvin` instruction. Increasing the array dimension results in larger blocks of memory requested per `mvin` over TileLink. Doubling the systolic array dimensions doubles the effective memory bandwidth and quadruples the compute throughput. Depending on how

much reuse there is within a layer and according to the tiling factors, the expected performance boost can be anywhere from 2 $\times$ -4 $\times$ .

For all the models, before feeding the data into the systolic array, the operands are zeropadded so that their dimensions are multiples of the size of the systolic array. In most of our benchmarks, this resulted in negligible added overhead of multiplying zeros. This overhead was highest in MobileNet, where it consumed 10% of the workload, but it significantly dropped with larger DNNs like Resnet.

We also found that due to their low arithmetic intensity and large memory footprint, depthwise convolutions would require feeding inputs sequentially into the systolic array, which would limit their performance. Therefore, we perform depthwise convolutions in MobileNet on the host processor itself. Prior work has demonstrated that the low arithmetic intensity of depthwise convolution can be an impediment to the efficient acceleration of MobileNet. This is also demonstrated in the results of our DSE - while depthwise convolution layers take up 18% of the runtime on our CPU implementation, they take up nearly 100% of the execution time in the accelerated

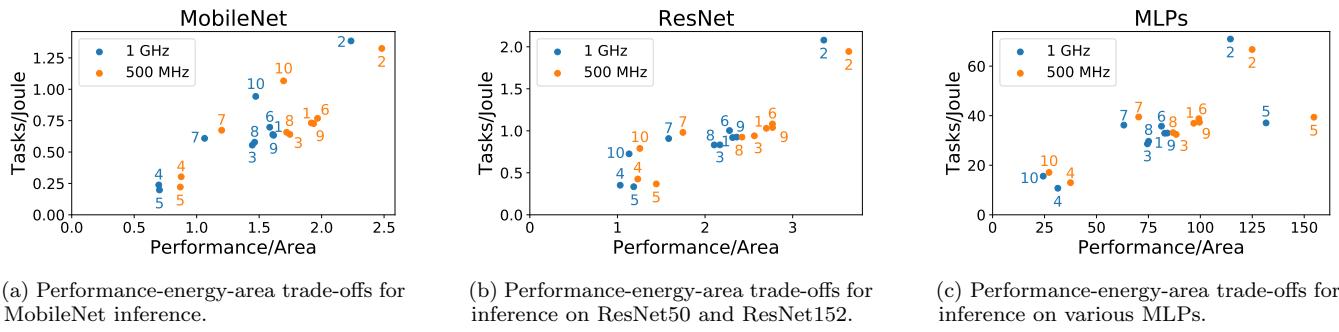


Figure 8: The tasks completed per unit energy as a function of performance per area for all our design points. The numeric labels correspond to rows in Table 1.

workload.

We observe from Figure 7b that the performance improvement between MLP topologies varies wildly, as a consequence of the shapes of their layers. The shapes of the layers affected not only the input/weight reuse, but also the amount by which GEMMs could be tiled in the scratchpad. The maximum tiling factors were a function of the scratchpad size and array dimensions, but narrow or non-divisible layers often reduced tiling factors, which also reduced performance. As an example, MLP 4 outperformed MLP 3, because its dimensions, which were powers-of-2 mapped better onto our maximum tiling factors.

### 3.4 Design Space Analysis

We integrate our results for power, performance, and area in Figure 8. We plot the performance per joule of each workload against the performance per unit area. In workloads which were not CPU-limited, such as MLPs, increasing memory capacity improved Gemmini’s area and energy efficiency by minimizing main memory accesses and improving locality within the scratchpad. In all workloads, design points which demanded extra memory bandwidth, such as ⑤ which increased bitwidth to 32 bits, reduced energy and area efficiency significantly.

Although the weight-stationary dataflow (②) did not noticeably improve performance, it did increase energy and area efficiency by removing the 32-bit accumulators within each PE, which greatly reduced the power consumption of the systolic mesh. The  $32 \times 32$  design (⑤), on the other hand, suffered from very low efficiency despite its high performance, because of its greatly increased power consumption and area overhead. The BOOM processor (⑩) improved energy- and area-efficiency significantly on MobileNet, which was severely CPU-limited, while it suffered on other workloads where the CPU was not the bottleneck.

Additionally, 500 MHz designs were generally more energy- and area-efficient than 1 GHz designs. The 500 MHz designs used more high-voltage-threshold (HVT) logic gates, reducing leakage power consumption more than enough to compensate for their slower compute performance.

## 4. PHYSICAL DESIGN

We perform an evaluation of the physical design properties of the generated accelerator design, to explore the area and power requirements of different design points, as well as to evaluate the place-and-route feasibility of Gemmini accelerators.

Since the end of Dennard scaling, power density has proven to be a significant factor in the design of digital systems-on-chip. In particular, custom accelerators with a dense collection of compute units (MACs, in the case of a matrix multiplication unit) are known to be sensitive to such thermal and energy constraints. As an example, the Google TPUv3 uses liquid cooling to assist the power dissipation from its dense array of compute units [33]. As such, it is important that the design space of such a matrix multiplication unit be evaluated through full physical design VLSI flows (placement and routing), allowing for the evaluation of the feasibility of the RTL design (timing closure), as well as power density and energy under different floorplans. Furthermore, physical design of the selected design points was necessary for the integration of the accelerator into test systems-on-chip for fabrication.

We evaluated four design points based on the results of the earlier DSE, by choosing two Gemmini configurations and place-and-routing each of them at both 500 MHz and 1 GHz. The first configuration is a  $16 \times 16$  systolic array, with dual dataflows, a 4-banked 256 KiB scratchpad memory and a 64 KiB accumulator. The second configuration is a  $32 \times 32$  systolic array, with dual dataflows, 4 banked 512 KiB scratchpad memory and 128 KiB accumulator.

Each of our selected design points was also evaluated using two different floorplans in TSMC 16 nm FinFET process technology. The first floorplan (Figure 9) organizes the accelerator’s SRAMs in a major block, leaving space on the side for the systolic mesh, as well as a routing channel across the block. The second floorplan, in Figure 10, organizes the accelerator’s SRAMs in a semi-ring around the computational mesh.

The placed designs are shown in Figures 9 and 10. In both floorplans, the controller was placed next to the Rocket host processor since the processor interacts with the controller to send instructions and data. Each

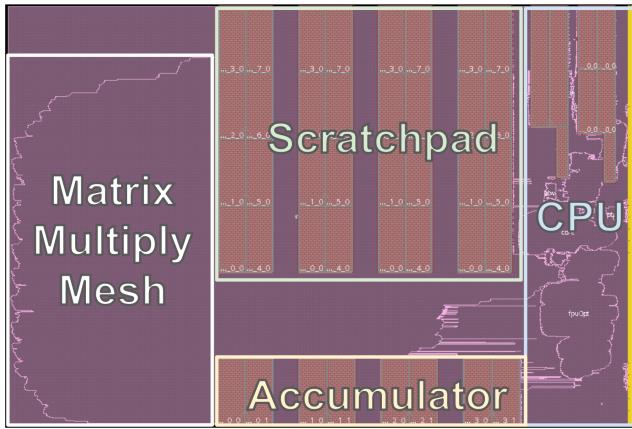


Figure 9: Block floorplan of the systolic array in TSMC 16nm for a  $16 \times 16$  design point. Rectangular blocks represent SRAM macros. Yellow points on the right are I/O pins for clock, reset, interrupt and external memory.

floorplan has intuitive benefits: while the block floorplan provides the systolic mesh more vertical access to the SRAM address lines, the semi-ring floorplan allows for more surface area contact between the systolic mesh and the SRAMs.

The comparison results are presented in Table 2. We can observe that the  $16 \times 16$  arrays performed similarly with both floorplans, even when increasing the frequency from 500 MHz to 1 GHz. For example, all  $16 \times 16$  design points achieved nearly the same Worst Negative Slack (WNS) for any given frequency, although the Total Negative Slack (TNS) was 79% higher with the semi-ring floorplan at 1 GHz. Additionally, the worst WNS for  $16 \times 16$  designs was -14 ps, which can easily be adjusted to meet timing requirements. The power consumption of  $16 \times 16$  designs also showed little variation across different floorplans, differing by only 1-3%. We can conclude from this that the generator is flexible enough to allow for a variety of floorplans, with reasonable a Quality-of-Result (QoR) for each.

However, we observe that for larger design points such as a  $32 \times 32$  configuration, the difference between the floorplans begins to have more noticeable impacts upon our timing results. For example, the semi-ring floorplan achieves a 41% better setup WNS at 1 GHz than the block floorplan, as well as an 81% better TNS, because wires between the scratchpad's SRAMs and the systolic mesh can be routed to travel a shorter average distance. However, neither of the  $32 \times 32$  floorplans were able to meet timing at 1 GHz, and the setup violations, which were several hundred picoseconds long, were significant enough to require several iterations of physical design attempts to have a chance of closing timing.

Overall, our physical design evaluation demonstrates that semi-ring floorplans can reduce power consumption and achieve faster clock frequencies. They do this by reducing wire lengths and by placing SRAMs where it is easier for physical design tools to route them to the systolic mesh. Furthermore, with  $16 \times 16$  designs, semi-

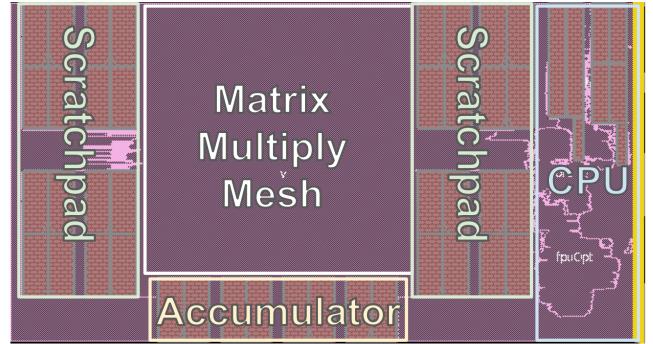


Figure 10: Semi-ring floorplan of the systolic array in TSMC 16nm for a  $16 \times 16$  design point.

ring floorplans also reduce area requirements. However, as the systolic array grows to  $32 \times 32$ , the area of the semi-ring floorplan grows faster than the block design. Thus, as an accelerator design grows, it may be necessary to change to a floorplan which clusters SRAMs closer together to meet area requirements.

## 5. DISCUSSION

Embedding and integrating the Gemmini generator within a composable, silicon-proven, open-source platform such as Rocket Chip allows for seamless integration with additional system components such as complex out-of-order cores and vector accelerators. We demonstrate integration with a Rocket RISC-V core, representative of low-energy processors found in embedded devices, as well as integration with BOOM, a high-performance OoO core, which revealed the power/performance tradeoffs when using a beefier core.

Gemmini was designed as a flexible generator for systolic GEMM accelerators to identify power, performance, and area trends as various parameters are varied, rather than to achieve state-of-the-art ML inference performance. Gemmini's DSE revealed the benefit of specializing the hardware for a weight-stationary-only dataflow as is the case in the TPU [6], and the system-level evaluation demonstrated that a larger scratchpad is only valuable if a workload is not CPU-limited.

Gemmini targets the most common kernel across many network architectures: matrix multiplications. A large portion of ML inference time is spent on fully connected layers [34], which are implemented as matrix multiplications. Furthermore, compute intensive convolutions used in CNNs can be efficiently mapped to matrix multiplications. By targeting GEMMs, Gemmini can to adapt to different network architectures and layer types, in contrast to specializing for convolutional layers. Additionally, GEMMs are a useful linear algebra primitive, which can expose the Gemmini generator to other application domains such as scientific computing.

Due to the speed limitations of RTL software simulation, some prior works choose to evaluate single-layer performance, and then extrapolate to report the performance of a full DNN. However, extrapolation of layer-by-layer performance neglects to consider shared-cache

Table 2: Floor-planned Design Points in TSMC 16nm FinFET Process Technology.

Design	Freq (MHz)	Floorplan	Area (mm <sup>2</sup> )	Power (mW)	Setup WNS (ps)	Setup TNS (ps)
16×16	500	Block	1.34	321.71	1	0
16×16	500	Semi-Ring	1.21	312.41	0	0
16×16	1000	Block	1.34	773.70	-12	-235
16×16	1000	Semi-Ring	1.21	766.12	-14	-420
32×32	500	Block	2.81	1058.24	-14	-100
32×32	500	Semi-Ring	3.01	1078.71	-9	-59
32×32	1000	Block	2.81	2796.51	-530	-2716
32×32	1000	Semi-Ring	3.01	2683.01	-315	-508

state between the host processor and the accelerator, as well as host-processor time between layers.

Furthermore, FPGA prototypes used to evaluate ASIC DNN accelerators often connect directly to an on-FPGA DRAM controller and thus see a higher memory throughput than an ASIC implementation would. This can distort performance and energy numbers. Gemmini has been evaluated using FireSim, which accurately models a last-level cache and DRAM to preserve simulation fidelity while executing on an FPGA.

Full-system simulation of RTL implementations is important not only for performance evaluation, but for functional validation as well. While an original version of the Gemmini design passed many individual benchmarks and micro-benchmarks, some design decisions such as methods for handling exceptions, memory management race conditions and TLB flushes were exposed only in a full-system multi-process environment such as Linux.

## 6. RELATED WORK

Systolic architectures first came into prominence in the early 1980s [9, 10], and since then, many systolic accelerators have been developed. There has also been much work on algorithms which can design new systolic arrays methodologically, rather than through ad-hoc intuition [35].

Early systolic arrays were used to compute convolutions [36], solutions to triangular linear systems [9], matrix multiplications, and more. Systolic architectures enable modular and extensible designs, use local neighbor-to-neighbor message passing, and contain easy-to-floorplan regular structures.

Systolic architectures have recently regained popularity, since the convolution and matrix multiplication kernels common in machine learning and deep learning application are highly susceptible to multi-dimensional acceleration using systolic arrays.

Commercially deployed ASIC implementations of NN accelerators include the Google TPU [6] for cloud workloads, as well as edge inference implementations by Samsung [7], Nvidia [37], Apple [38], and Tesla [8, 39]. In particular, a detailed description of the original TPU

implementation includes a  $256 \times 256$  matrix multiplication unit implemented using a reduced-precision systolic MAC array with a weight stationary dataflow for NN inference in the cloud. Successor versions included floating-point representation, additional memory, and improved utilization for both training and inference [33].

Prior work has demonstrated the integration of an open-source commercial DNN accelerator (NVDLA) with the Rocket Chip ecosystem and the FireSim platform [40]. The accelerator in this work was integrated using the memory bus, as opposed to Gemmini which is integrated using the RoCC interface. Prior work [41] has also demonstrated the integration of academic NN accelerators with the Rocket Chip ecosystem using the RoCC interface, but did not use systolic architectures for that purpose. Gemmini puts an emphasis on enabling design space exploration rather than single design-point integration.

Academic researchers have proposed numerous systolic accelerators, especially for neural-network inference. For example, NeuFlow [42] was a systolic-inspired architecture which allowed individual processing elements (PEs) to be re-configured at runtime to perform tasks such as multiply-accumulates, divisions, and non-linear activations. ShiDianNao [43], similarly, allowed PEs to be re-configured at runtime to perform multiply-accumulates, additions, and max poolings. Eyeriss [44] implemented a weight-stationary dataflow using a spatial array. Eyeriss v2 [45] improved on the original Eyeriss by demonstrating a new PE architecture that can operate on sparse CSC-encoded matrices, and a hierarchical mesh NoC capable of unicast, multicast, and broadcast data transfers to maximize reuse. These and other systolic-inspired architectures typically permit both global and local connections between PEs and global memory, which is not strictly systolic, but often improves performance.

Several previous proposals [13, 46, 47] have presented performance and energy benefits resulting from flexible data-flow options in NN accelerators. However, the benefits and impact of the dataflow structure of NN accelerators is still an active area of research, and some works [48] have shown that optimal memory-hierarchies

and loop-blocking strategies can have a more significant impact on energy efficiency than the choice of dataflows.

Various energy efficient neural network accelerator proposals have also been presented in the integrated circuits community [49–59]. Many of these proposals focus on exploiting sparsity and quantization features of DNNs. Furthermore, while some of these proposals address runtime-configurability, they still address only a single fabrication-capable design point, and most do not present design and elaboration time parameterization. Further, most of these accelerators are tested in isolation, often without a fully integrated software environment, hence potentially neglecting system-level effects.

A host of DNN accelerators targeted for FPGA implementation have also been proposed [60–71], taking advantage of FPGA reconfigurability to implement exotic layers, specialize the hardware for a specific network, and evaluate multiple design points. However, FPGA acceleration frameworks do not necessarily translate well to ASIC implementations, and are not ideal for scenarios where energy efficiency is critical.

Some prior works [13, 72–78] use analytical or high-level model-based simulations to evaluate different parameterizations of a proposed accelerator architecture. In contrast, Gemmini performs design space exploration on the RTL directly and uses feedback from FPGA-accelerated simulation and physical design to find optimal design points for ASIC implementation.

Since the energy consumed during DNN inference and matrix multiplication is often dominated by external memory accesses, academic researchers have proposed processing in memory [79–85]. These works include the development of new SRAM circuits and the use of novel devices such as ReRAMs. Gemmini is designed and validated for CMOS implementation, and uses design space exploration to discover the ideal memory access patterns and memory hierarchy to conserve energy.

Researchers have also proposed methodological systems and algorithms to automatically generate systolic architectures directly from the algorithms they are meant to accelerate. For example, PolySA [86] analyzes polyhedral models to attempt to find the optimal mapping between a sequential algorithm and a set of parallel PEs. Yang et al. [48] extended the Halide programming language to automatically generate C++ high-level-synthesis (HLS) implementations of systolic arrays.

Prior work has also introduced TVM [87] and VTA [25] as an integrated research platform for SW/HW evaluation of NN accelerators. While Gemmini and VTA hold many architectural similarities, including the use of a GEMM core, explicit memory management, and explicit instruction dependency handling, VTA has primarily targeted FPGA accelerators implementations, as opposed to Gemmini which currently targets primarily ASIC designs and has been used in the fabrication on multiple test-chips. Furthermore, Gemmini’s integration with the RISC-V eco-system enables an additional level of customization in SW/HW co-design.

## 7. FUTURE WORK

The Gemmini generator has been used in the fabrication of two test system-on-chips. The chips were taped-out within approximately a month of each other in different process technologies, demonstrating the flexibility and utility of the Gemmini generator. Further evaluation of the integration of the Gemmini accelerators within the context of these larger embedded vision processors will be performed when the chips complete the fabrication process.

As demonstrated previously, CPU operations can significantly slow down inference on our workloads. Operations which map convolutions to GEMMs, such as `im2col`, can make up a significant portion of this overhead. To address this we intend to map convolutions to GEMMs transparently in hardware.

Some additional overheads come from zero-padding matrices so that their dimensions can tile onto the systolic array, which reduces utilization at the boundaries of our arrays. By breaking up a single, large, systolic array into numerous smaller ones operating in parallel, we can possibly reduce zero-padding requirements while still preserving the same compute throughput [88].

Finally, a generator-based methodology can be useful for the hardware/software co-design process through the integration of hardware generator and compiler parameters. Future integration with optimizing DSLs and compilers such as Halide or TVM will allow for better code generation which considers the generator parameters, hence allowing for better cross-layer data re-use and optimization.

## 8. CONCLUSION

This work presented Gemmini, an open source and agile systolic array generator that enables systematic evaluations of deep-learning architectures. This systematic evaluation is demonstrated through a DSE case study, identifying bottlenecks in common DNN inference workloads, and capturing the variation of performance improvements of different workloads running on different hardware configurations. With a baseline design equipped with a  $16 \times 16$  systolic array, Gemmini demonstrated  $90\times$  and  $70\times$  inference speedups on ResNet-152 and ResNet-50, respectively, when compared to a cache-optimized CPU implementation, and two to three orders of magnitude speedup on MLPs. We demonstrate the critical importance of full-system evaluation by showing that even though an accelerator can effectively accelerate individual layers of DNNs, it often fails to achieve impressive performance improvements on the entire DNN if any part of it is not efficiently mapped onto the accelerator. For example, although a Gemmini baseline design was able to accelerate the first layer of MobileNet by  $330\times$ , it failed to accelerate the entire network beyond  $6\times$  using a Rocket host processor and  $18\times$  using a BOOM host processor, due to the presence of depthwise convolutions. We also show that even with DNNs that have similar network architectures, performance may vary based upon the shape and size of different layers. Looking forward, we believe Gemmini will enable a new range of systematic evaluations and HW/SW co-design

of deep learning workloads.

## 9. REFERENCES

- [1] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *nature*, vol. 521, no. 7553, p. 436, 2015.
- [2] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in neural information processing systems*, pp. 1097–1105, 2012.
- [3] I. Sutskever, O. Vinyals, and Q. V. Le, “Sequence to sequence learning with neural networks,” in *Advances in neural information processing systems*, pp. 3104–3112, 2014.
- [4] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, “Openai gym,” *arXiv preprint arXiv:1606.01540*, 2016.
- [5] J. Mahler, J. Liang, S. Niyaz, M. Laskey, R. Doan, X. Liu, J. A. Ojea, and K. Goldberg, “Dex-net 2.0: Deep learning to plan robust grasps with synthetic point clouds and analytic grasp metrics,” *arXiv preprint arXiv:1703.09312*, 2017.
- [6] N. P. Jouppi, C. Young, N. Patil, D. Patterson, G. Agrawal, R. Bajwa, S. Bates, S. Bhatia, N. Boden, A. Borchers, R. Boyle, P.-I. Cantin, C. Chao, C. Clark, J. Coriell, M. Daley, M. Dau, J. Dean, B. Gelb, T. V. Ghaemmaghami, R. Gottipati, W. Gulland, R. Hagnmann, C. R. Ho, D. Hogberg, J. Hu, R. Hundt, D. Hurt, J. Ibarz, A. Jaffey, A. Jaworski, A. Kaplan, H. Khaitan, D. Killebrew, A. Koch, N. Kumar, S. Lacy, J. Laudon, J. Law, D. Le, C. Leary, Z. Liu, K. Lucke, A. Lundin, G. MacKean, A. Maggiore, M. Mahony, K. Miller, R. Nagarajan, R. Narayanaswami, R. Ni, K. Nix, T. Norrie, M. Omernick, N. Penukonda, A. Phelps, J. Ross, M. Ross, A. Salek, E. Samadiani, C. Severn, G. Sizikov, M. Snelham, J. Souter, D. Steinberg, A. Swing, M. Tan, G. Thorson, B. Tian, H. Toma, E. Tuttle, V. Vasudevan, R. Walter, W. Wang, E. Wilcox, and D. H. Yoon, “In-datacenter performance analysis of a tensor processing unit,” in *Proceedings of the 44th Annual International Symposium on Computer Architecture, ISCA ’17*, (New York, NY, USA), pp. 1–12, ACM, 2017.
- [7] J. Song, Y. Cho, J.-S. Park, J.-W. Jang, S. Lee, J.-H. Song, J.-G. Lee, and I. Kang, “7.1 an 11.5 tops/w 1024-mac butterfly structure dual-core sparsity-aware neural processing unit in 8nm flagship mobile soc,” in *2019 IEEE International Solid-State Circuits Conference-(ISSCC)*, pp. 130–132, IEEE, 2019.
- [8] P. J. Bannon, K. A. Hurd, and E. Talpes, “Accelerated mathematical engine,” Jan. 24 2019. US Patent App. 15/710,433.
- [9] H.-T. Kung, “Why systolic architectures?,” *IEEE computer*, vol. 15, no. 1, pp. 37–46, 1982.
- [10] H. Kung and C. E. Leiserson, “Systolic arrays (for vlsi),” in *Sparse Matrix Proceedings 1978*, vol. 1, pp. 256–282, Society for Industrial and Applied Mathematics, 1979.
- [11] B. Nikolic, “Simpler, more efficient design,” in *ESSCIRC Conference 2015 - 41st European Solid-State Circuits Conference (ESSCIRC)*, pp. 20–25, Sep. 2015.
- [12] B. Nikolic, E. Alon, and K. Asanovic, “Generating the next wave of custom silicon,” in *ESSCIRC 2018 - IEEE 44th European Solid State Circuits Conference (ESSCIRC)*, pp. 6–11, Sep. 2018.
- [13] K. Kwon, A. Amid, A. Gholami, B. Wu, K. Asanovic, and K. Keutzer, “Co-design of deep neural nets and neural net accelerators for embedded vision applications,” in *2018 55th ACM/ESDA/IEEE Design Automation Conference (DAC)*, pp. 1–6, IEEE, 2018.
- [14] K. Asanovic, R. Avizienis, J. Bachrach, S. Beamer, D. Biancolin, C. Celio, H. Cook, D. Dabbel, J. Hauser, A. Izraelevitz, S. Karandikar, B. Keller, D. Kim, J. Koenig, Y. Lee, E. Love, M. Maas, A. Magyar, H. Mao, M. Moreto, A. Ou, D. A. Patterson, B. Richards, C. Schmidt, S. Twigg, H. Vo, and A. Waterman, “The rocket chip generator,” Tech. Rep. UCB/EECS-2016-17, EECS Department, University of California, Berkeley, Apr 2016.
- [15] C. Celio, D. A. Patterson, and K. Asanovic, “The berkeley out-of-order machine (boom): An industry-competitive, synthesizable, parameterized risc-v processor,” *EECS Department, University of California, Berkeley, Tech. Rep. UCB/EECS-2015-167*, 2015.
- [16] S. Karandikar, H. Mao, D. Kim, D. Biancolin, A. Amid, D. Lee, N. Pemberton, E. Amaro, C. Schmidt, A. Chopra, Q. Huang, K. Kovacs, B. Nikolic, R. Katz, J. Bachrach, and K. Asanovic, “Firesim: Fpga-accelerated cycle-exact scale-out system simulation in the public cloud,” in *2018 ACM/IEEE 45th Annual International Symposium on Computer Architecture (ISCA)*, pp. 29–42, June 2018.
- [17] J. Bachrach, H. Vo, B. Richards, Y. Lee, A. Waterman, R. Avizienis, J. Wawrzynek, and K. Asanović, “Chisel: Constructing hardware in a scala embedded language,” in *Proceedings of the 49th Annual Design Automation Conference, DAC ’12*, (New York, NY, USA), pp. 1216–1225, ACM, 2012.
- [18] B. Jacob, S. Kligys, B. Chen, M. Zhu, M. Tang, A. G. Howard, H. Adam, and D. Kalenichenko, “Quantization and training of neural networks for efficient integer-arithmetic-only inference,” *CoRR*, vol. abs/1712.05877, 2017.
- [19] H. M. Cook, A. S. Waterman, and Y. Lee, “Sifive tilelink specification,” tech. rep., SiFive Inc., 2018. <https://www.sifive.com/documentation/tilelink/tilelink-spec/>.
- [20] S. Han, H. Mao, and W. J. Dally, “Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding,” *arXiv preprint arXiv:1510.00149*, 2015.
- [21] Y. Lin, S. Han, H. Mao, Y. Wang, and W. J. Dally, “Deep gradient compression: Reducing the communication bandwidth for distributed training,” *arXiv preprint arXiv:1712.01887*, 2017.
- [22] Y. S. Shao, S. L. Xi, V. Srinivasan, G.-Y. Wei, and D. Brooks, “Co-designing accelerators and soc interfaces using gem5-aladdin,” in *2016 49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, pp. 1–12, IEEE, 2016.
- [23] I. Y. John L. Gustafson, “Beating floating point at its own game: Posit arithmetic,” *Supercomputing Frontiers and Innovations*, vol. 4, jun 2017.
- [24] J. E. Smith, “Decoupled access/execute computer architectures,” in *Proceedings of the 9th Annual Symposium on Computer Architecture, ISCA ’82*, (Los Alamitos, CA, USA), pp. 112–119, IEEE Computer Society Press, 1982.
- [25] T. Moreau, T. Chen, Z. Jiang, L. Ceze, C. Guestrin, and A. Krishnamurthy, “VTA: an open hardware-software stack for deep learning,” *CoRR*, vol. abs/1807.04188, 2018.
- [26] D. Biancolin, S. Karandikar, D. Kim, J. Koenig, A. Waterman, J. Bachrach, and K. Asanović, “Fased: Fpga-accelerated simulation and evaluation of dram,” in *The 2019 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays (FPGA’19)*, FPGA ’19, (New York, NY, USA), ACM, 2019.
- [27] D. Claudio Ciresan, U. Meier, L. M. Gambardella, and J. Schmidhuber, “Deep big simple neural nets excel on handwritten digit recognition,” *arXiv preprint arXiv:1003.0358*, 2010.
- [28] U. Meier, D. C. Ciresan, L. M. Gambardella, and J. Schmidhuber, “Better digit recognition with a committee of simple neural nets,” in *2011 International Conference on Document Analysis and Recognition*, pp. 1250–1254, IEEE, 2011.
- [29] X. Lu, Y. Tsao, S. Matsuda, and C. Hori, “Speech enhancement based on deep denoising autoencoder.,” in *Interspeech*, pp. 436–440, 2013.
- [30] J. Ngiam, A. Khosla, M. Kim, J. Nam, H. Lee, and A. Y. Ng, “Multimodal deep learning,” in *Proceedings of the 28th international conference on machine learning (ICML-11)*, pp. 689–696, 2011.
- [31] S. Hadjis, F. Abuzaid, C. Zhang, and C. Ré, “Caffe con troll: Shallow ideas to speed up deep learning,” in *Proceedings of*

- the Fourth Workshop on Data Analytics in the Cloud*, DanaC'15, (New York, NY, USA), pp. 2:1–2:4, ACM, 2015.
- [32] S. Chetlur, C. Woolley, P. Vandermersch, J. Cohen, J. Tran, B. Catanzaro, and E. Shelhamer, “cudnn: Efficient primitives for deep learning,” *CoRR*, vol. abs/1410.0759, 2014.
- [33] Google, “System architecture | cloud tpu | google cloud.” <https://cloud.google.com/tpu/docs/system-architecture>, 2019. Accessed: 2019-07-15.
- [34] J. Park, M. Naumov, P. Basu, S. Deng, A. Kalaiah, D. Khudia, J. Law, P. Malani, A. Malevich, S. Nadathur, J. Pino, M. Schatz, A. Sidorov, V. Sivakumar, A. Tulloch, X. Wang, Y. Wu, H. Yuen, U. Diril, D. Dzhulgakov, K. Hazelwood, B. Jia, Y. Jia, L. Qiao, V. Rao, N. Rotem, S. Yoo, and M. Smelyanskiy, “Deep learning inference in facebook data centers: Characterization, performance optimizations and hardware implications,” 2018.
- [35] A. Iványi, ed., *Algorithms of Informatics*, vol. 2. Kempelen Farkas Hallgatói Információk Központ, 2011.
- [36] H. Kung and S. W. Song, “A Systolic 2-D Convolution Chip,” tech. rep., CARNEGIE-MELLON UNIV PITTSBURGH PA DEPT OF COMPUTER SCIENCE, 1981.
- [37] NVIDIA, “Nvidia turing gpu architecture.” <https://www.nvidia.com/content/dam/en-zz/Solutions/design-visualization/technologies/turing-architecture/NVIDIA-Turing-Architecture-Whitepaper.pdf>, 2018. Accessed: 2019-08-01.
- [38] Apple, “The future is here: iphone x.” <https://www.apple.com/newsroom/2017/09/the-future-is-here-iphone-x/>, 2017. Accessed: 2019-08-01.
- [39] P. J. Bannon and K. A. Hurd, “Systems and methods for hardware-based pooling,” July 4 2019. US Patent App. 15/862,369.
- [40] F. Farshchi, Q. Huang, and H. Yun, “Integrating nvidia deep learning accelerator (nvdla) with risc-v soc on firesim,” *CoRR*, 2019.
- [41] S. Eldridge, A. Waterland, M. Seltzer, J. Appavoo, and A. Joshi, “Towards general-purpose neural network computing,” in *2015 International Conference on Parallel Architecture and Compilation (PACT)*, pp. 99–112, Oct 2015.
- [42] C. Farabet, B. Martini, B. Corda, P. Akselrod, E. Culurciello, and Y. LeCun, “Neuflow: A runtime reconfigurable dataflow processor for vision,” in *CVPR Workshops*, pp. 109–116, 2011.
- [43] Z. Du, R. Fasthuber, T. Chen, P. Ienne, L. Li, T. Luo, X. Feng, Y. Chen, and O. Temam, “Shidiannao: Shifting vision processing closer to the sensor,” in *2015 ACM/IEEE 42nd Annual International Symposium on Computer Architecture (ISCA)*, pp. 92–104, June 2015.
- [44] Y.-H. Chen, T. Krishna, J. S. Emer, and V. Sze, “Eyeriss: An energy-efficient reconfigurable accelerator for deep convolutional neural networks,” *IEEE Journal of Solid-State Circuits*, vol. 52, no. 1, pp. 127–138, 2017.
- [45] Y. Chen, T. Yang, J. Emer, and V. Sze, “Eyeriss v2: A flexible accelerator for emerging deep neural networks on mobile devices,” *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, vol. 9, pp. 292–308, June 2019.
- [46] W. Lu, G. Yan, J. Li, S. Gong, Y. Han, and X. Li, “Flexflow: A flexible dataflow accelerator architecture for convolutional neural networks,” in *2017 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, pp. 553–564, Feb 2017.
- [47] F. Tu, S. Yin, P. Ouyang, S. Tang, L. Liu, and S. Wei, “Deep convolutional neural network architecture with reconfigurable computation patterns,” *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 25, pp. 2220–2233, Aug 2017.
- [48] X. Yang, M. Gao, J. Pu, A. Nayak, Q. Liu, S. Bell, J. Setter, K. Cao, H. Ha, C. Kozyrakis, and M. Horowitz, “DNN dataflow choice is overrated,” *CoRR*, vol. abs/1809.04070, 2018.
- [49] K. Ueyoshi, K. Ando, K. Hirose, S. Takamaeda-Yamazaki, J. Kadomoto, T. Miyata, M. Hamada, T. Kuroda, and M. Motomura, “Quest: A 7.49tops multi-purpose log-quantized dnn inference engine stacked on 96mb 3d sram using inductive-coupling technology in 40nm cmos,” in *2018 IEEE International Solid - State Circuits Conference - (ISSCC)*, pp. 216–218, Feb 2018.
- [50] J. Lee, C. Kim, S. Kang, D. Shin, S. Kim, and H. Yoo, “Unpu: A 50.6tops/w unified deep neural network accelerator with 1b-to-16b fully-variable weight bit-precision,” in *2018 IEEE International Solid - State Circuits Conference - (ISSCC)*, pp. 218–220, Feb 2018.
- [51] D. Bankman, L. Yang, B. Moons, M. Verhelst, and B. Murmann, “An always-on 3.8pj/86% cifar-10 mixed-signal binary cnn processor with all memory on chip in 28nm cmos,” in *2018 IEEE International Solid - State Circuits Conference - (ISSCC)*, pp. 222–224, Feb 2018.
- [52] T. Karnik, D. Kurian, P. Aseron, R. Dorrance, E. Alpman, A. Nicoara, R. Popov, L. Azarenkov, M. Moiseev, L. Zhao, S. Ghosh, R. Misoczki, A. Gupta, M. Akhila, S. Muthukumar, S. Bhandari, Y. Satish, K. Jain, R. Flory, C. Kanthanapanit, E. Quijano, B. Jackson, H. Luo, S. Kim, V. Vaidya, A. Elsherbini, R. Liu, F. Sheikh, O. Tickoo, I. Klotchkov, M. Sastry, S. Sun, M. Bhartiya, A. Srinivasan, Y. Hoskote, H. Wang, and V. De, “A cm-scale self-powered intelligent and secure iot edge mote featuring an ultra-low-power soc in 14nm tri-gate cmos,” in *2018 IEEE International Solid - State Circuits Conference - (ISSCC)*, pp. 46–48, Feb 2018.
- [53] D. Shin, J. Lee, J. Lee, and H. Yoo, “14.2 dnpu: An 8.1tops/w reconfigurable cnn-rnn processor for general-purpose deep neural networks,” in *2017 IEEE International Solid-State Circuits Conference (ISSCC)*, pp. 240–241, Feb 2017.
- [54] S. Yin, P. Ouyang, S. Tang, F. Tu, X. Li, L. Liu, and S. Wei, “A 1.06-to-5.09 tops/w reconfigurable hybrid-neural-network processor for deep learning applications,” in *2017 Symposium on VLSI Circuits*, pp. C26–C27, June 2017.
- [55] K. Ando, K. Ueyoshi, K. Orimo, H. Yonekawa, S. Sato, H. Nakahara, M. Ikebe, T. Asai, S. Takamaeda-Yamazaki, T. Kuroda, and M. Motomura, “Breibin memory: A 13-layer 4.2 k neuron/0.8 m synapse binary/ternary reconfigurable in-memory deep neural network accelerator in 65 nm cmos,” in *2017 Symposium on VLSI Circuits*, pp. C24–C25, June 2017.
- [56] C. Kim, S. Kang, D. Shin, S. Choi, Y. Kim, and H. Yoo, “A 2.1tflops/w mobile deep rl accelerator with transposable pe array and experience compression,” in *2019 IEEE International Solid- State Circuits Conference - (ISSCC)*, pp. 136–138, Feb 2019.
- [57] A. Sayal, S. Fathima, S. S. T. Nibhanupudi, and J. P. Kulkarni, “14.4 all-digital time-domain cnn engine using bidirectional memory delay lines for energy-efficient edge computing,” in *2019 IEEE International Solid- State Circuits Conference - (ISSCC)*, pp. 228–230, Feb 2019.
- [58] J. Lee, J. Lee, D. Han, J. Lee, G. Park, and H. Yoo, “7.7 lnpu: A 25.3tflops/w sparse deep-neural-network learning processor with fine-grained mixed precision of fp8-fp16,” in *2019 IEEE International Solid- State Circuits Conference - (ISSCC)*, pp. 142–144, Feb 2019.
- [59] J. Yue, R. Liu, W. Sun, Z. Yuan, Z. Wang, Y. Tu, Y. Chen, A. Ren, Y. Wang, M. Chang, X. Li, H. Yang, and Y. Liu, “A 65nm 0.39-to-140.3tops/w 1-to-12b unified neural network processor using block-circulant-enabled transpose-domain acceleration with 8.1 × higher tops/mm<sup>2</sup>and 6t/hbst-tram-based 2d data-reuse architecture,” in *2019 IEEE International Solid- State Circuits Conference - (ISSCC)*, pp. 138–140, Feb 2019.
- [60] H. Zeng, R. Chen, C. Zhang, and V. Prasanna, “A framework for generating high throughput CNN implementations on FPGAs,” in *Proceedings of the 2018 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays - FPGA '18*, ACM Press, 2018.
- [61] Y. Wang, J. Xu, Y. Han, H. Li, and X. Li, “DeepBurning: automatic generation of FPGA-based learning accelerators

- for the neural network family,” in *Proceedings of the 53rd Annual Design Automation Conference on - DAC ’16*, ACM Press, 2016.
- [62] J. Shen, Y. Huang, Z. Wang, Y. Qiao, M. Wen, and C. Zhang, “Towards a uniform template-based architecture for accelerating 2d and 3d CNNs on FPGA,” in *Proceedings of the 2018 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays - FPGA ’18*, ACM Press, 2018.
- [63] R. DiCecco, G. Lacey, J. Vasiljevic, P. Chow, G. Taylor, and S. Areibi, “Caffeinated FPGAs: FPGA framework for convolutional neural networks,” in *2016 International Conference on Field-Programmable Technology (FPT)*, IEEE, dec 2016.
- [64] Y. Guan, H. Liang, N. Xu, W. Wang, S. Shi, X. Chen, G. Sun, W. Zhang, and J. Cong, “FP-DNN: An automated framework for mapping deep neural networks onto FPGAs with RTL-HLS hybrid templates,” in *2017 IEEE 25th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*, IEEE, apr 2017.
- [65] K. Guo, L. Sui, J. Qiu, J. Yu, J. Wang, S. Yao, S. Han, Y. Wang, and H. Yang, “Angel-eye: A complete design flow for mapping CNN onto embedded FPGA,” *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 37, pp. 35–47, jan 2018.
- [66] X. Zhang, J. Wang, C. Zhu, Y. Lin, J. Xiong, W.-m. Hwu, and D. Chen, “Dnnbuilder: An automated tool for building high-performance dnn hardware accelerators for fpgas,” in *Proceedings of the International Conference on Computer-Aided Design, ICCAD ’18*, (New York, NY, USA), pp. 56:1–56:8, ACM, 2018.
- [67] C. Zhang, P. Li, G. Sun, Y. Guan, B. Xiao, and J. Cong, “Optimizing fpga-based accelerator design for deep convolutional neural networks,” in *Proceedings of the 2015 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays, FPGA ’15*, (New York, NY, USA), pp. 161–170, ACM, 2015.
- [68] S. I. Venieris and C. Bouganis, “fpgaconvnet: A framework for mapping convolutional neural networks on fpgas,” in *2016 IEEE 24th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*, pp. 40–47, May 2016.
- [69] H. Sharma, J. Park, D. Mahajan, E. Amaro, J. K. Kim, C. Shao, A. Mishra, and H. Esmaeilzadeh, “From high-level deep neural models to fpgas,” in *Microarchitecture (MICRO), 2016 49th Annual IEEE/ACM International Symposium on*, pp. 1–12, IEEE, 2016.
- [70] J. Zhang and J. Li, “Improving the performance of opencl-based fpga accelerator for convolutional neural network,” in *ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, 2017.
- [71] C. Wang, L. Gong, Q. Yu, X. Li, Y. Xie, and X. Zhou, “Dlau: A scalable deep learning accelerator unit on fpga,” *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 36, pp. 513–517, March 2017.
- [72] R. Yazdani, A. Segura, J. Arnau, and A. Gonzalez, “An ultra low-power hardware accelerator for automatic speech recognition,” in *2016 49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, pp. 1–12, Oct 2016.
- [73] L. Song, X. Qian, H. Li, and Y. Chen, “Pipelayer: A pipelined reram-based accelerator for deep learning,” in *2017 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, pp. 541–552, Feb 2017.
- [74] P. Srivastava, M. Kang, S. K. Gonugondla, S. Lim, J. Choi, V. Adve, N. S. Kim, and N. Shanbhag, “Promise: An end-to-end design of a programmable mixed-signal accelerator for machine-learning algorithms,” in *2018 ACM/IEEE 45th Annual International Symposium on Computer Architecture (ISCA)*, pp. 43–56, June 2018.
- [75] S. Sharify, A. D. Lascorz, K. Siu, P. Judd, and A. Moshovos, “Loom: Exploiting weight and activation precisions to accelerate convolutional neural networks,” in *Proceedings of the 55th Annual Design Automation Conference, DAC ’18*, (New York, NY, USA), pp. 20:1–20:6, ACM, 2018.
- [76] S. Angizi, Z. He, and D. Fan, “Dima: A depthwise cnn in-memory accelerator,” in *2018 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*, pp. 1–8, Nov 2018.
- [77] C. Min, J. Mao, H. Li, and Y. Chen, “Neuralhmc: An efficient hmc-based accelerator for deep neural networks,” in *Proceedings of the 24th Asia and South Pacific Design Automation Conference, ASPDAC ’19*, (New York, NY, USA), pp. 394–399, ACM, 2019.
- [78] T. Nowatzki, V. Gangadhar, N. Ardalani, and K. Sankaralingam, “Stream-dataflow acceleration,” in *2017 ACM/IEEE 44th Annual International Symposium on Computer Architecture (ISCA)*, pp. 416–429, June 2017.
- [79] J. Liu, H. Zhao, M. A. Ogleari, D. Li, and J. Zhao, “Processing-in-memory for energy-efficient neural network training: A heterogeneous approach,” in *2018 51st Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, pp. 655–668, Oct 2018.
- [80] P. Chi, S. Li, C. Xu, T. Zhang, J. Zhao, Y. Liu, Y. Wang, and Y. Xie, “Prime: A novel processing-in-memory architecture for neural network computation in reram-based main memory,” in *2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture (ISCA)*, pp. 27–39, June 2016.
- [81] H. Yan, H. R. Cherian, E. C. Ahn, and L. Duan, “Celia: A device and architecture co-design framework for stt-mram-based deep learning acceleration,” in *Proceedings of the 2018 International Conference on Supercomputing, ICS ’18*, (New York, NY, USA), pp. 149–159, ACM, 2018.
- [82] H. Yan, A. H. Aboutalebi, and L. Duan, “Efficient allocation and heterogeneous composition of nvm crossbar arrays for deep learning acceleration,” in *2018 IEEE 37th International Performance Computing and Communications Conference (IPCCC)*, pp. 1–8, Nov 2018.
- [83] Y. Zha, E. Nowak, and J. Li, “Liquid silicon: A nonvolatile fully programmable processing-in-memory processor with monolithically integrated reram for big data/machine learning applications,” in *2019 Symposium on VLSI Circuits*, 2019.
- [84] Y. Ji, Y. Zhang, X. Xie, S. Li, P. Wang, X. Hu, Y. Zhang, and Y. Xie, “Fpsa: A full system stack solution for reconfigurable reram-based nn accelerator architecture,” in *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS ’19*, (New York, NY, USA), pp. 733–747, ACM, 2019.
- [85] L. Chang, X. Ma, Z. Wang, Y. Zhang, W. Zhao, and Y. Xie, “Corn: In-buffer computing for binary neural network,” in *2019 Design, Automation Test in Europe Conference Exhibition (DATE)*, pp. 384–389, March 2019.
- [86] J. Cong and J. Wang, “Polysa: polyhedral-based systolic array auto-compilation,” in *2018 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*, pp. 1–8, IEEE, 2018.
- [87] T. Chen, T. Moreau, Z. Jiang, L. Zheng, E. Yan, H. Shen, M. Cowan, L. Wang, Y. Hu, L. Ceze, C. Guestrin, and A. Krishnamurthy, “TVM: An automated end-to-end optimizing compiler for deep learning,” in *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*, (Carlsbad, CA), pp. 578–594, USENIX Association, 2018.
- [88] H. T. Kung, “Don’t use a single large systolic array, use many small ones instead.” Workshop on ML for Systems at ISCA, 2019.